

SIGMOD Officers, Committees, and Awardees

Chair

Raghu Ramakrishnan
Yahoo! Research
2821 Mission College
Santa Clara, CA 95054
USA
<First8CharsOfLastName AT
yahoo-inc.com>

Vice-Chair

Yannis Ioannidis
University of Athens
Department of Informatics & Telecom
Panepistimioupolis, Informatics Buildings
157 84 Ilissia, Athens
HELLAS
<yannis AT di.uoa.gr>

Secretary/Treasurer

Mary Fernández
ATT Labs - Research
180 Park Ave., Bldg 103, E277
Florham Park, NJ 07932-0971
USA
<mff AT research.att.com>

SIGMOD Executive Committee:

Curtis Dyreson, Mary Fernández, Yannis Ioannidis, Phokion Kolaitis, Alexandros Labrinidis, Lisa Singh, Tamer Özsu, Raghu Ramakrishnan, and Jeffrey Xu Yu.

Advisory Board: Tamer Özsu (Chair), University of Waterloo, <tozsu AT cs.uwaterloo.ca>, Rakesh Agrawal, Phil Bernstein, Peter Buneman, David DeWitt, Hector Garcia-Molina, Jim Gray, Masaru Kitsuregawa, Jiawei Han, Alberto Laender, Krithi Ramamritham, Hans-Jörg Schek, Rick Snodgrass, and Gerhard Weikum.

Information Director:

Jeffrey Xu Yu, The Chinese University of Hong Kong, <yu AT se.cuhk.edu.hk>

Associate Information Directors:

Marcelo Arenas, Denilson Barbosa, Ugur Cetintemel, Manfred Jeusfeld, Alexandros Labrinidis, Dongwon Lee, Michael Ley, Rachel Pottinger, Altigran Soares da Silva, and Jun Yang.

SIGMOD Record Editor:

Alexandros Labrinidis, University of Pittsburgh, <labrinid AT cs.pitt.edu>

SIGMOD Record Associate Editors:

Magdalena Balazinska, Denilson Barbosa, Ugur Çetintemel, Brian Cooper, Andrew Eisenberg, Cesar Galindo-Legaria, Leonid Libkin, Jim Melton, Len Seligman, and Marianne Winslett.

SIGMOD DiSC Editor:

Curtis Dyreson, Washington State University, <cdyreson AT eecs.wsu.edu>

SIGMOD Anthology Editor:

Curtis Dyreson, Washington State University, <cdyreson AT eecs.wsu.edu>

SIGMOD Conference Coordinators:

Lisa Singh, Georgetown University, <singh AT cs.georgetown.edu>

PODS Executive: Phokion Kolaitis (Chair), IBM Almaden, <kolaitis AT almaden.ibm.com>, Foto Afrati, Catriel Beeri, Georg Gottlob, Leonid Libkin, and Jan Van Den Bussche.

Sister Society Liaisons:

Raghu Ramakrishnan (SIGKDD), Yannis Ioannidis (EDBT Endowment).

Awards Committee: Gerhard Weikum (Chair), Max-Planck Institute of Computer Science, <weikum AT mpi-sb.mpg.de>, Peter Buneman, Mike Carey, David Maier, and Moshe Y. Vardi.

SIGMOD Officers, Committees, and Awardees (continued)

SIGMOD Edgar F. Codd Innovations Award

For innovative and highly significant contributions of enduring value to the development, understanding, or use of database systems and databases. Until 2003, this award was known as the "SIGMOD Innovations Award." In 2004, SIGMOD, with the unanimous approval of ACM Council, decided to rename the award to honor Dr. E.F. (Ted) Codd (1923 - 2003) who invented the relational data model and was responsible for the significant development of the database field as a scientific discipline. Recipients of the award are the following:

Michael Stonebraker (1992)	Jim Gray (1993)	Philip Bernstein (1994)
David DeWitt (1995)	C. Mohan (1996)	David Maier (1997)
Serge Abiteboul (1998)	Hector Garcia-Molina (1999)	Rakesh Agrawal (2000)
Rudolf Bayer (2001)	Patricia Selinger (2002)	Don Chamberlin (2003)
Ronald Fagin (2004)	Michael Carey (2005)	Jeffrey D. Ullman (2006)
Jennifer Widom (2007)		

SIGMOD Contributions Award

For significant contributions to the field of database systems through research funding, education, and professional services. Recipients of the award are the following:

Maria Zemankova (1992)	Gio Wiederhold (1995)	Yahiko Kambayashi (1995)
Jeffrey Ullman (1996)	Avi Silberschatz (1997)	Won Kim (1998)
Raghu Ramakrishnan (1999)	Michael Carey (2000)	Laura Haas (2000)
Daniel Rosenkrantz (2001)	Richard Snodgrass (2002)	Michael Ley (2003)
Surajit Chaudhuri (2004)	Hongjun Lu (2005)	Tamer Özsu (2006)
Hans-Jörg Schek (2007)		

SIGMOD Doctoral Dissertation Award

The annual ACM SIGMOD Doctoral Dissertation Award, inaugurated in 2006, recognizes excellent research by doctoral candidates in the database field.

- **2006 Winner:** Gerome Miklau, University of Washington
Runners-up: Marcelo Arenas, University of Toronto; Yanlei Diao, University of California at Berkeley.
- **2007 Winner:** Boon Thau Loo, University of California at Berkeley
Honorable Mentions: Xifeng Yan, University of Illinois at Urbana-Champaign; Martin Theobald, Saarland University

A complete listing of all SIGMOD Awards is available at: <http://www.sigmod.org/awards/>

[Last updated on March 12, 2007]

Editor's Notes

Welcome to the December 2007 issue of SIGMOD Record. We start this issue with a welcome message by Yannis Ioannidis, aimed primarily at new SIGMOD members (many of whom joined SIGMOD right after the SIGMOD/PODS 2007 conference). Following Yannis' message, you will find a short note by Curtis Dyreson about the SIGMOD Anthology volume that you are receiving together with this issue.

Next, we have two regular **articles**, which are both "*critiques*". The first one, is on physical database design in general and the TAB benchmark in particular (by Nicolas Bruno). The second article is on nulls, three-valued logic and ambiguity in SQL (by Claude Rubinson).

We continue with an article in the **Surveys Column** (edited by Cesar Galindo-Legaria), on *Context Models* (by Bolchini, Curino, Quintarelli, Schreiber, and Tanca). With context-aware systems becoming increasingly pervasive in everyday life, this data-oriented survey should be an excellent starting point on the topic.

Next we have an article on the **Systems and Prototypes Column** (edited by Magdalena Balazinska), about the *Intel Mash Maker*, which was one of the demos in SIGMOD 2007. The previous Systems and Prototypes article was published in the September 2004 issue of SIGMOD Record. I am very happy to see the column revitalized again and feature descriptions of exciting and innovative systems and prototypes.

The **Distinguished Profiles in Data Management Column** (edited by Marianne Winslett) features an interview of Ricardo Baeza-Yates who is Vice President of Yahoo! Research in Europe and Latin America. Read Ricardo Baeza-Yates' interview to find out (among many other things) about CS Research in Latin America, his multi-continent commute for Yahoo! and how to get real data in academia.

We continue with an article in the **Research Centers Column** (edited by Ugur Cetintemel), about *Data and Web Management Research at Politecnico di Milano* (by the 16 members of the research group, which include Stefano Ceri, Cristiana Bolchini, Piero Fraternali, Fabio A. Schreiber, and Letizia Tanca). The article highlights the group's research across two different dimensions: *data-driven research* (which includes work on context-aware and mobile databases; the survey on context models in this issue is authored by members of this group) and *web-driven research*.

Next we have three articles in the **Event Reports Column** (edited by Brian Cooper). The first is the *Report on the First International Workshop on Ranking in Databases (DBRank'07)* which was held in April 2007, together with ICDE 2007. The second is the *Report on the Fourth International Workshop on Data Management for Sensor Networks (DMSN 2007)*, which was held in September 2007, together with VLDB 2007. The third is the *Report on the first VLDB workshop on Management of Uncertain Data (MUD)*, which was also held in September 2007, together with VLDB 2007.

We also have two important **announcements** in this issue: the Call for Nominations for the 2008 ACM SIGMOD Awards (deadline: April 7, 2008), and the Call for Submissions for the SIGMOD 2008 Undergraduate Research Poster Competition (deadline: April 4, 2008).

We continue with a very important **Call for Participation** in the Tribute to Honor Jim Gray, which will be held on May 31, 2008 at UC Berkeley.

With the December issue being late, we were able to include for the first time the **Calls for Papers** for almost all the workshops that will be held together with this year's SIGMOD conference. These are (in alphabetical order):

- DaMoN 2008: 4th International Workshop on Data Management on New Hardware (deadline: April 11, 2008),
- DBTest 2008: 1st International Workshop on Testing Database Systems (deadline: April 11, 2008),
- MobiDE 2008: 7th International ACM Workshop on Data Engineering for Wireless and Mobile Access (deadline: March 26, 2008),
- WebDB 2008: 11th International Workshop on the Web and Databases (deadline: April 6, 2008),
- XIME-P 2008: 5th International Workshop on XQuery Implementation, Experience and Perspectives (deadline: March 28, 2008).

When I had in mind that the June 2007 issue would be a collectors' item, I only thought this would be the case because of the change of the cover design and the change in Editor. Little did I know that there would be a third reason. Unfortunately, there was a problem with one of the papers' math fonts during printing (specifically: they were not included) which made the paper incomprehensible. We are reprinting the paper (*Estimating the Selectivity of tf-idf based Cosine Similarity Predicates*, by Tata and Patel) in its entirety in this issue (under the **Errata Column**), with the math included this time. Additionally, ACM and the printer have been able to "debug" this problem and have made all the necessary arrangements for this not to happen again in the future.

Alexandros Labrinidis
February 2007

Welcome Message to New SIGMOD Members

On behalf of the entire Executive Committee of ACM SIGMOD, it is a great pleasure for me to be writing to you for the first time since you have become members. As the Vice-Chair of this Special Interest Group (SIG) and responsible for members' issues, I will be in touch with you at regular intervals, informing you of any major developments, important activities, and new initiatives we may be undertaking. This will be happening by direct email on a periodic basis (roughly every three months) as well as through the quarterly issues of SIGMOD Record; at the same time, our website (www.sigmod.org) will always be up to date with the latest information.

Your benefits are all analyzed in detail on the SIGMOD website. Regarding the content you receive as part of them, our general philosophy is for trying to provide as much as possible online, over the web, and avoid the expensive and environment-unfriendly media and paper printing and shipping. Current copyright difficulties, however, prevent us from making everything available in this fashion; hence, several levels of membership have been established that provide content on different media. We are in a continuous effort to strike agreements with all relevant copyright owners to increase what is provided online and we will be informing you on any developments on that front. Likewise, I would like to be hearing from you (MyFirstName @ di.uoa.gr) regarding any comments, ideas, or thoughts you might have about how we may improve your benefits as SIGMOD members.

SIGMOD is one of the largest SIGs within ACM. By joining it, you are becoming part of a group of over 2400 scientists from all across the globe, whose goal is to promote research and technological advancement in the field of data, information, and knowledge management. The future calls for a continuous move towards the so-called "information / knowledge societies"; SIGMOD has a critical role to play in facilitating and promoting the development of the appropriate technologies that would realize the positive aspects of such societies and protect against their negative aspects. This can only be achieved if all of us join forces and collectively steer relevant activities in fruitful directions, so your active involvement in SIGMOD is important. Especially if you are a student, we need your fresh ideas and hope for your regular participation in the SIGMOD-sponsored conferences and workshops.

I want to welcome you again to ACM SIGMOD and look forward to hearing from you and to seeing you at SIGMOD/PODS 2008 in Vancouver.

Sincerely,
Yannis Ioannidis
SIGMOD Vice-Chair

SIGMOD Anthology Volume 6

Curtis Dyreson
ACM SIGMOD Anthology Editor
Utah State University

We are pleased to include a DVD of Volume 6 of the ACM SIGMOD Anthology with this issue of SIGMOD Record. The Anthology is a digital library for the database research community developed by ACM SIGMOD with past cooperation of the VLDB Endowment, the IEEE Technical Committee on Data Engineering, and the EDBT Endowment, and with the current assistance of many publishers and individuals, both in providing permission to include material in the Anthology, and in locating physical copies to scan into digital form. The bibliographic information in the Anthology is integrated with the DBLP Bibliography.

Volume 6 of the Anthology includes the following.

- *ACM Transactions on Information Systems* (1983-2005)
- *ACM Transactions on Internet Technology* (2001-2005)
- Proceedings of the Australasian Database Conference (2002-2005)
- Proceedings of the Asia-Pacific Conference on Conceptual Modelling (2004-2005)
- Reports of various CODASYL Committees from the Charles Babbage Institute
- An interview with Charles W. Bachman
- Books by Richard T. Snodgrass and Gio Wiederhold
- The DBLP Browser and Bibliography

I wish to thank Michael Ley, the Anthology's founding editor, for helping to prepare this volume, and for integrating the Anthology with DBLP, which is a wonderful resource for the research community. I also wish to thank the special contributors to this volume: Gio Wiederhold, Richard T. Snodgrass, the members of the CODASYL Systems and DBTG Committees (in particular Charles W. Bachman and T. William Olle), Elisabeth Kaplan and Carrie Seib of the Charles Babbage Institute (Center for the History of Information Technology) at the University of Minnesota, and the ACM.

A Critical Look at the TAB Benchmark for Physical Design Tools

Nicolas Bruno

Microsoft Research

nicolasb@microsoft.com

Abstract

There has recently been considerable research on physical design tuning algorithms. At the same time, there is only one published methodology to evaluate the quality of different, competing approaches: the TAB benchmark. In this paper we describe our experiences with TAB. We first report an experimental evaluation of TAB on our latest prototype for physical design tuning. We then identify certain weakness in the benchmark and briefly comment on alternatives to improve its usefulness.

1 Introduction

Lately there has been considerable effort in the database community on reducing the total cost of ownership of database installations. Specifically, physical design tuning has become relevant, and most vendors nowadays include automated tools to tune database physical designs as part of their products (e.g., [3, 10, 14]). Given a query workload W and a storage budget B , these tools find the set of physical structures (or configuration) that fits in B and results in the lowest cost for W (see Figure 1).

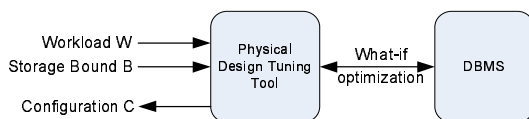


Figure 1: Architecture of Physical Design Tools.

Although there has been considerable research in new algorithms to find good configurations and extensions to newer physical structures (e.g., [2, 4, 5, 7, 8, 11, 12, 15]), much less attention has been paid on methodologies to evaluate the quality of different approaches. Specifically, we are aware of only one publication that proposes a benchmark of physical design tools: the Toronto Automatic Benchmark, or TAB for short [9]. In this paper we describe our experiences with TAB when evaluating the quality of different alternatives, both in the context of a shipping product [3] and also on different experimental prototypes that we implemented over the last few years [5, 6]. Specifically, in Section 2 we review the TAB benchmark [9]. In Section 3 we report an experimental

evaluation of TAB. Then, in Section 4 we analyze both the results of the experimental evaluation and also the benchmark itself. In doing so, we identify certain weaknesses in the design of TAB (specifically, on the benchmark metrics, the choice of baseline configurations, and some combination of database/workloads) and briefly comment on alternatives to mitigate their impact.

2 The TAB Benchmark

Reference [9] introduces a framework to evaluate the quality of automated physical design tuners, which we refer to as *TAB*. We next review the three components of the benchmark: the evaluation metrics, a baseline configuration to compare against recommendations, and the set of databases/workloads to tune.

Evaluation Metric Consider a workload W over a database D , and suppose that a tuner recommends configuration C for W . *TAB* evaluates the quality of C using $\mathcal{M}_{C,W}$, which returns, for an input time t , the number of queries in W that executed faster than t :

$$\mathcal{M}_{C,W}(t) = \frac{|\{q \in W : cost(q, C) \leq t\}|}{|W|}$$

where $cost(q, C)$ is the actual execution time of query q under configuration C . For pragmatic purposes, a timeout T_{max} is chosen and $cost(q, C)$ is capped by T_{max} . Therefore, it is always $\mathcal{M}_{C,W}(T_{max}) = 1$.

Baseline Configuration *TAB* identifies a special configuration, called $1C$, which consists of all single-column indexes over the database tables. Reference [9] justifies the choice of $1C$ by stating that “... the consistently good performance of the single column configuration suggests a practical improvement of DBMS configuration recommends...”, “... $1C$ was also far better than the configurations recommended by both systems...”, and “...a very conservative overall workload assessment results in $1C$ producing almost 17 times better results than $R!$ ”.

Databases and Workloads *TAB* uses two databases. The first one is a publicly available non-redundant reference protein database [13], or *NREF* for short, which

provides a collection of protein sequence data from several genome sequencing projects. The second one is the TPC-H benchmark used to evaluate the performance of database systems [1]. The workloads in [9] are chosen to “...represent fragments of typical iceberg queries, that is, queries that compute aggregate functions over a set of attributes to find aggregate values satisfying certain conditions, grouped in different ways”. A typical query for the reference protein database is shown below¹:

```
SELECT T1.nref_id, COUNT(DISTINCT T2.nref_id)
FROM taxonomy T1, taxonomy T2, protein P
WHERE T1.taxon_id = T2.taxon_id AND
      T1.nref_id = P.nref_id AND
      P.p_name = 'Phosphotransferase'
GROUP BY T1.nref_id
```

For the TPC-H database, *TAB* does not use the QGen workload of [1], but rather one that mimics that of *NREF*.

3 Running *TAB*

We now report an experimental evaluation of *TAB* in our physical design tuning prototype based on [5]. Our objective with this experiment was two-fold. First, we wanted to analyze the performance of our prototype design tuner and compare its quality with the baseline configuration of [9]. Second, we wanted to understand the design decisions behind *TAB*, and question whether there was room for improvement in the benchmark definition itself. We used a Intel Xeon 3.2 GHz CPU with 2GB of RAM (we allocated 1GB of RAM to the DBMS for the experiments) and a 250GB, 7200rpm hard drive to store data. We used Microsoft SQL Server 2005 as the database engine. For each workload, we proceeded as follows. Following [9], we first created three copies of the original database, and deployed a different configuration on each instance. The first one, which we denote by adding a suffix *-P* to the database name, has only primary indexes. The second one, which we denote by adding a suffix *-PIC* to the database name, additionally contains all valid single-column indexes². The third one, which we denote by adding a suffix *-R* to the database name, is obtained by running our physical design tuning tool for the input workload considering both clustered and non-clustered indexes with a storage bound equal to the size of the *-PIC* configuration. Table 1 shows statistics on the databases and workloads. (Note that we also evaluated the original TPC-H workload generated using the QGen utility.)

To avoid external factors in skewing the results, we performed the following additional steps. First, we stopped

¹All queries in the workload follow the same pattern: (i) self-join of a table T_1 , (ii) join with a table T_2 that has a selective predicate, (iii) aggregates on the values of T_1 tables.

²Restrictions in the DBMS prevent us from creating certain indexes, such as indexes with keys larger than 900 bytes.

Database	Size	# Indexes
NREF-P	8GB	6
NREF-PIC	34GB	35
NREF-R (tuned with NREF3J)	28GB	31
TPC-H-P	12GB	8
TPC-H-PIC	34GB	61
TPC-H-R (tuned with UnTH3J)	21GB	29
TPC-H-R (tuned with QGen[1])	34GB	15

Table 1: Databases used in the evaluation.

all non-essential operating system services to avoid interference. Second, we defragmented both the disk where data resided and also the indexes inside the database. Third, we created the same set of statistics in all databases. Finally, we executed each query five times –with cold buffers– and kept the median execution time. We used a timeout T_{max} of 30 minutes as in [9], but no execution exceeded T_{max} .

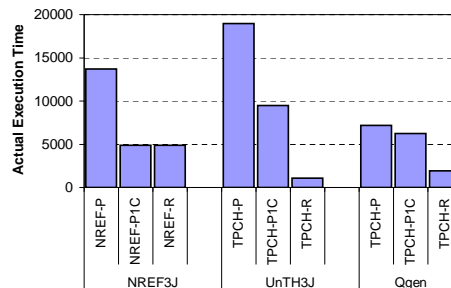


Figure 2: Overall Execution Times.

Figure 2 shows the overall execution times for all workloads and databases. Figure 3 shows $\mathcal{M}(t)$ for each database/workload combination. Finally, Figure 4 shows a variation of the \mathcal{M} metric where we used the optimizer’s estimated cost rather than the actual execution cost for the queries in the workload. We analyze these results next.

4 Analyzing *TAB*

We now analyze the results of the experimental evaluation of the previous section. In doing so, we also address some issues on *TAB* that we found during the evaluation and comment on some alternatives to diminish their impact.

Overall Comments. Figure 2 shows that the recommended configurations resulted in substantial improvement over the basic *-P* configurations. Specifically, the improvements were 64% for *NREF/NREF3J*, 94% for *TPC-H/UnTH3J*, and 73% for *TPC-H/QGen*. A noticeable difference with [9] is the performance of the *-PIC* configurations. While for *NREF/NREF3J* both *-PIC* and *-R* resulted in roughly the same performance, for *TPC-H/UnTH3J* the performance of *-PIC* lies almost exactly between that of *-P* and *-R*. Also, for *TPC-H/QGen* the performance of *-PIC*

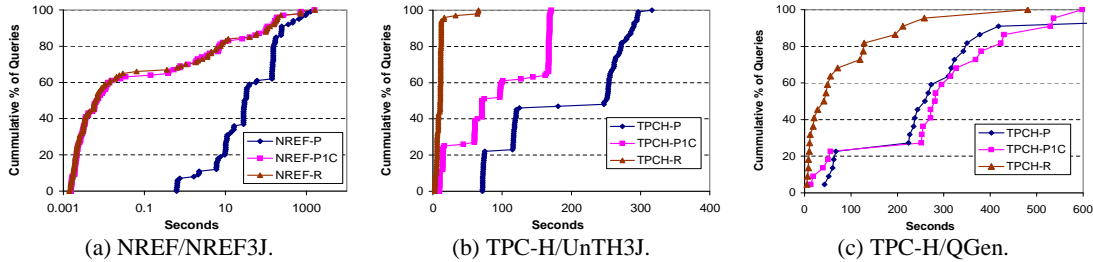


Figure 3: Actual execution times for varying databases and workloads.

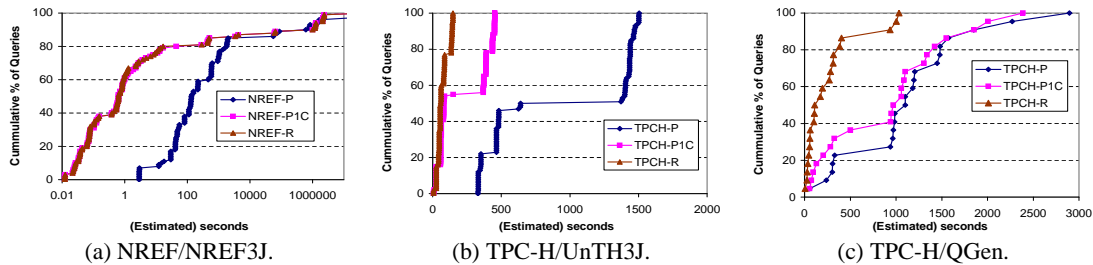


Figure 4: Optimizer estimated execution times for varying databases and workloads.

is only slightly better than that of *-P* (a 13% improvement compared with 73% of *-R*). Figures 3 and 4 give additional information about the relative performance of different configurations. Almost 80% of the queries in *NREF-3J* finished in less than 10 seconds under either *-PIC* or *-R*, but only 10% of the queries under *-P* finished in that amount of time. For *TPC-H/UnTH3J*, all 100% of the queries finished in 75 seconds or less under *R*, where the percentages were 50% for *-PIC* and only 5% for *-P*. Finally, for *TPC-H/QGen*, 90% of the queries ran in less than 220 seconds under *-R*, where only 22% of the queries did the same under either *-P* or *-PIC*. Interestingly, the \mathcal{M} curves for *-P* and *-PIC* cross each other for *TPC-H/QGen* in Figure 3(c), and therefore it is not clear how to interpret their relative performance beyond our original claim that the configurations were comparable. We examine and comment on the design of the TAB benchmark itself next.

4.1 Evaluation Metrics

The metric used to compare tuners is a crucial component of a benchmark. Usually, the existing literature uses a single number to measure the quality of recommendations, called *percentage improvement*, and defined as $1 - \text{actual cost} / \text{recommended cost}$. TAB recognizes that a single number might not provide enough detail to thoroughly evaluate and compare physical design tuners, and proposes the \mathcal{M} metric to address this limitation. While we agree with the deficiency pointed out in [9] regarding single-value metrics, we identify some problems in \mathcal{M} .

4.1.1 Actual vs. Estimated Cost

The \mathcal{M} metric is based on the actual time it takes to execute queries in the workload. We believe that in the con-

text of evaluating a full system (i.e., not only the tuning tool, but also the query optimizer, query processor, and even the underlying operating system) this is clearly the best, most unbiased choice. However, if the purpose is an isolated evaluation of physical design tools, we claim that execution costs are, although important, less relevant. The reason is that using execution costs potentially introduces additional variables that are outside the scope of the evaluated tool. We next clarify this claim with real examples.

The Role of the Optimizer. It is important to note the we are bound to execute what the optimizer decides it is the best plan for a given query³. Consider the following example, simplified from a real query in *NREF/NREF3J*:

```
SELECT R.* FROM R, S
WHERE predicate(R) AND R.x=S.y
```

and suppose that the optimizer estimates that only a handful of tuples from *R* satisfy `predicate(R)`. If an index on *S.y* is available, the optimizer would find that a nested-index-loop alternative that first gets all valid tuples from *R* and then fetches the matches from *S* might be a better alternative than, say, a hash join. Now suppose that the estimate is not right due to limitations in the optimizer's cost model, and in reality almost all tuples in *R* satisfy `predicate(R)`. In this case, the index-nested-loop plan, although it is costed the lowest by the optimizer and therefore chosen if possible, would execute much slower than the sub-optimal (to the eyes of the optimizer) hash-join alternative. Now the problem is clear. Consider the query above under the *-P* and *-PIC* configurations. The optimizer would pick the hash-join based alternative under *-P* (because there is no index on *S.y* in *-P*) and the

³Hints can be used to override optimizer's decisions, but should be used with caution and as a last resource

index-nested-loop alternative under *-PIC* (because the index is present). The net effect is that the execution cost under *-PIC* would be significantly larger than that under *-P*, and we would tend to rank the tuner that produced *-P* higher than the one that produced *-PIC*. However, note that under *-PIC* the optimizer *considered* the hash-join alternative but discarded it in favor of the index-nested-loop plan! In fact, within the optimizer’s cost model, the index-nested-loop alternative is better than the hash-based alternative in both *-PIC* and *-P* (although the former plan is not implementable under *-P*).

When purely evaluating the *quality of a physical design tuner*, we should be careful to freeze any external variables. It is therefore reasonable to assume that the optimizer is correct and the physical design tool exploits accurate information. Using the optimizer’s expected cost rather than the actual execution cost of queries has precisely that effect, provided that the optimizer is operating under the same statistical model for all configurations (which we can achieve by materializing the same set of statistics, including those that are associated with indexes, in each database instance).

Runtime conditions. Another problem when using actual execution times is the unwanted presence of external factors that can compromise the accuracy of the measurements. In one of our earlier experiments, we noticed that the execution cost of a plan under *-P* was twice as fast as the corresponding plan under *-R* (which was odd since *-R* contained a strict superset of the indexes in *-P* and the query did not do any updates). Even more puzzling, a closer inspection of both plans revealed that they were indeed identical. After a long debugging session, we realized that the root cause of the problem was index fragmentation. In fact, the query required a sequential scan over an index. Since the index under *-P* was not fragmented, the execution engine could go through the index using sequential I/O, which is fast. In contrast, under *-R* the execution engine had to do one random I/O every 5 disk blocks on average due to fragmentation in the index, which resulted in a larger execution time overall.

It seems unfair to punish a tuner tool due to external factors that are not under its control. Although we minimized this effect by defragmenting the indexes and underlying disk in our experiments, there is always a chance that external factors play a role in biasing the results.

4.1.2 Timeouts in the \mathcal{M} Metric

Reference [9] introduces a timeout value T_{max} that caps the maximum execution time of a query, set as 30 minutes. Although this is a practical issue to avoid very long running queries, it introduces some problems in the benchmark methodology. Specifically, it changes *a-posteriori* the optimization function that has been agreed upon and

leveraged in tuning tools. Consider the following extreme scenario, with a 2-query workload that contains a light query q_1 , which executes in 5 seconds under *-P* and a heavy query q_2 that executes in 3,600 seconds under *-P*. Consider a tuner T_1 that optimizes q_2 as much as possible at the expense of not fully optimizing q_1 , and suppose that the resulting times are $(q_1=4, q_2=1900)$, with an overall execution time of 1,905 seconds, or a 47% improvement. A second tuning tool T_2 , knowing *in advance* the 1,800-second timeout value, might optimize q_1 without considering q_2 obtaining the following times $(q_1=1, q_2=3600)$, with an overall execution time of 3,601 seconds, or just 0.1% improvement. Considering timeouts, the results are $(q_1=4, q_2=T_{max})$ for T_1 vs. $(q_1=1, q_2=T_{max})$ for T_2 , harshly underestimating T_1 ’s quality.

We believe that timeouts open the door for the possibility of “cheating” the benchmark by tools that exploit the subtle issues described above, and therefore recommend against using timeouts when evaluating configurations. (Strictly speaking, \mathcal{M} itself uses a different optimization criterium to what has been adopted in tuning tools, but its limitations are less severe than those derived from timeout values.)

4.1.3 Aggregating individual results

Once we obtain execution times for each query in the workload, we need to display this information in a meaningful manner. TAB therefore introduced the \mathcal{M} metric to show detailed information about performance of physical tuners. This metric is interesting in the sense that (i) allows to compare multiple tuners simultaneously, and (ii) allows for certain goal-oriented evaluation (such as 30% of the queries should execute in sub-second time [9]). One drawback of the \mathcal{M} metric is that it does not report per-query comparisons because the individual queries are sorted in different orders. It is not possible, just by looking at \mathcal{M} to draw conclusions about the performance of specific queries. For instance, although some queries were better under *-P* than under *-PIC* for *NREF*, Figure 3(a) is not enough to show this fact.

We next propose a complementary metric, which we call \mathcal{I} , that focuses on query-by-query performance. Consider configurations C_1 and C_2 coming from two tuning tools. We then compute, for each query q_i in the workload, the value $v_i = cost(q_i, C_1) - cost(q_i, C_2)$. Clearly, positive v_i values correspond to queries that were better under C_1 than under C_2 , and negative v_i values correspond to the opposite situation. We then sort v_i values and plot the results. Figures 5(a-c) show our proposed metric for the databases/workloads in our evaluation. Analogously, Figures 5(d-f) shows a variation of the \mathcal{I} metric that normalizes each v_i value by $cost(q_i, -P)$ (i.e., the cost of the query under the configuration that only has primary indexes). We can quickly see, for instance, that for

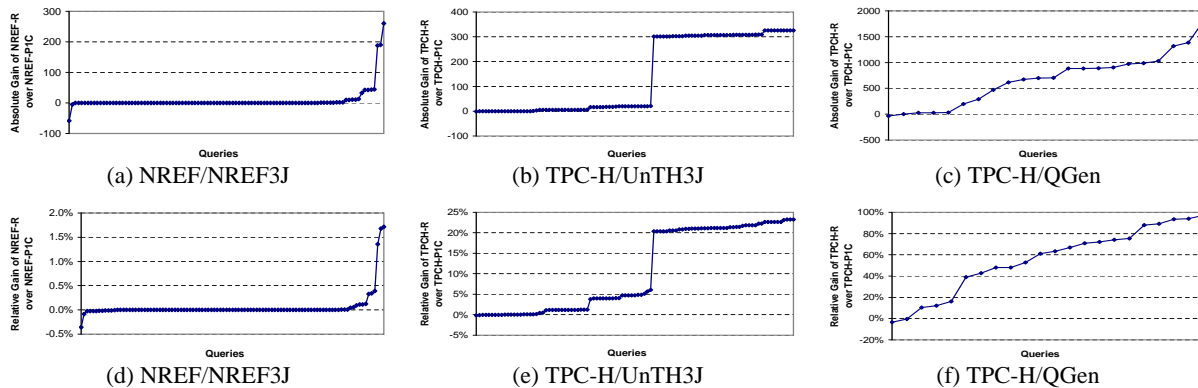


Figure 5: Proposed \mathcal{I} metric to compare physical design tuners.

NREF/NREF3J both *-PIC* and *-R* result in almost no difference in performance, but there are still some queries (which are easily identified in the figure) for which *-R* resulted in better performance. Also, for *TPC-H/UnTH3J* we can see that there are two clusters of queries: one that results in almost no variation between *-PIC* and *-R*, and another for which the variation is significant in *-R*'s favor. Finally, *TPC-H/QGen* goes from no variation to almost 100% relative change in performance.

Although the \mathcal{I} metric gives additional information on a per-query basis, it cannot be used to compare more than two configurations. We believe that \mathcal{M} and \mathcal{I} are complementary metrics that provide different types of insights when comparing physical design tuners.

4.2 Baseline Configuration

Before beginning our experiments we were surprised by the consistently good performance of *-PIC* claimed in [9]. Our experiments led to two key observations. First, current tuning tools result in configurations that range from comparable to *-PIC* to significantly better than *-PIC*. Second, there is a very large variance of performance of *-PIC* configurations, ranging from close to the best known solutions to close to the trivial configurations. In light of these observations, and based on Figures 3 and 4, we argue against using *-PIC* as a baseline configuration to compare against recommendations.

At some level, it is intuitive that *-PIC* would not be particularly helpful in general, and specifically for decision support workloads that require aggregating or filtering multiple columns. However, *-PIC* is essentially indistinguishable from the best recommended configuration for the *NREF/NREF3J* instance, which features queries with joins and aggregation. We next explain the main reasons behind this rather unexpected result.

Implied Index Columns. Secondary indexes in a DBMS store at the leaf nodes enough information to locate tuples in the primary index. To avoid storing record-

ids, which are volatile in the presence of updates, modern systems use the columns in the primary index as this identifier⁴. This implies that, for all practical purposes, single-column indexes implicitly behave as multi-column indexes. We cannot seek these implied columns, but execution plans can rely on them as if they were explicitly declared. Now consider the *NREF* database. Not only the tables in *NREF* are narrow (the median number of columns is only five), but also the primary indexes are wide. As an example, consider table *source*, which is composed of six columns, four of which are part of the primary index. In this case, every single-column index on *source* essentially contains 4 or 5 out of the 6 columns of the table! In fact, since just a minority of the table columns is not present in the index leaf nodes, single-column indexes in *-PIC* actually behave like “covering-indexes” for *NREF*.

Workload. Even for the “quasi”-covering-indexes in *-PIC* there are very simple examples that result in bad execution plans. Consider the following query in *NREF*:

```
SELECT taxon_id_2
FROM neighboring_seq
WHERE nref_id_2 < 'NF00000300'
```

where the predicate filters all but 7531 rows. The recommended configuration for this query has a covering index on (*nref_id_2*, *taxon_id*), so it can seek the relevant tuples and return the results optimally with an expected time of 0.51 units and an actual execution time of 0.078 seconds. Note that the primary index for table *neighboring_seq* does not contain column *taxon_id*. Therefore, *-PIC* cannot use the index on *nref_id_2* to locate the valid tuples and then fetch the remaining columns because the cost would be too high. Instead, the best plan for *-PIC* is to scan the index on *taxon_id*, which implicitly contains column *nref_id_2* and filter on the fly the resulting tuples. The expected cost of this strategy is 3.22 units (632 times slower than *-R*), and the actual execution

⁴If the primary index is not unique, a special “uniquifier” column is implicitly added.

time is 67.6 seconds (8667 times slower than *-R*). Additionally, for workloads with many updates, the performance of *-PIC* would be heavily deteriorated due to the overhead of updating the relevant indexes. Clearly, *-PIC* can result in very bad execution plans for the simplest of queries. A closer analysis of *NREF3J* shows, however, that for virtually all queries such situations fortunately do not happen, and thus *-PIC* performs extremely well in this scenario.

4.3 Database/Workloads

Once the metrics have been defined, the most important component of a benchmark is the actual databases and workloads over which it would be run. The TAB benchmark goes in the right direction by proposing both real (*NREF*) and synthetic (*TPC-H*) databases and workloads. However, it is also an example of how careful we need to be when designing benchmarks: by only considering *NREF/NREF3J* and *TPC-H/UnTH3J*, reference [9] arrives at the questionable conclusion that *-PIC* is a very competitive configuration. Another subtle problem with the *NREF* workload is that there is over six orders of magnitude difference between the slowest and fastest queries. Having very long queries in the workload is that these “rogue” queries might bias the result, specially in conjunction with timeout values in the \mathcal{M} metric.

We believe that database/workload generation for the purposes of physical design benchmarks is an open area of research. In the meantime, we believe that useful benchmarks should contain databases/workloads taken from at least the following three “buckets”:

- Micro-benchmarks that evaluate the different capabilities of the underlying DBMS and for which optimal configurations can be manually derived.
- Synthetic, complex workloads that exercise the full capabilities of the underlying query processor and cannot be manually analyzed.
- Real databases and workloads to address subtle scenarios that might have been overlooked in the previous two buckets.

5 Conclusions

In this paper we reported an experimental evaluation of the TAB benchmark for automated physical design tuners. We described TAB and its design choices and analyzed the quality of recommendations of our prototypes for the databases and workloads specified in TAB. In doing so, we identified certain weaknesses in the design of TAB and proposed alternatives to mitigate their impact. While TAB is an important first step in the area of physical tuning tool benchmarking, we believe that more work is needed. In particular, one of the biggest challenges in the area is to

obtain a principled way to generate databases and workloads that are comprehensive enough to compare competing tools that might be based on very different principles. We note that both [9] and this work assume that the underlying database system does not change across alternative physical design tuners. If this assumption does not hold, it is not even clear how the different tuners could/should be compared (actual execution times might be an ultimate metric, but they evaluate the whole system rather than just the tuning tool). We believe that this is a rather deep problem that might have profound implications in future research on physical design tuning.

References

- [1] TPC Benchmark H. Available at <http://www.tpc.org>.
- [2] S. Agrawal, S. Chaudhuri, and V. Narasayya. Automated selection of materialized views and indexes in SQL databases. In *Proceedings of VLDB*, 2000.
- [3] S. Agrawal et al. Database Tuning Advisor for Microsoft SQL Server 2005. In *Proceedings of VLDB*, 2004.
- [4] S. Agrawal, V. Narasayya, and B. Yang. Integrating vertical and horizontal partitioning into automated physical database design. In *Proceedings of SIGMOD*, 2004.
- [5] N. Bruno and S. Chaudhuri. Automatic physical database tuning: A relaxation-based approach. In *Proceedings of SIGMOD*, 2005.
- [6] N. Bruno and S. Chaudhuri. To tune or not to tune? A Lightweight Physical Design Alerter. In *Proceedings of VLDB*, 2006.
- [7] S. Chaudhuri, M. Datar, and V. Narasayya. Index selection for databases: A hardness study and a principled heuristic solution. In *IEEE Trans. Knowl. Data Eng.* 16(11), 2004.
- [8] S. Chaudhuri and V. Narasayya. An efficient cost-driven index selection tool for Microsoft SQL Server. In *Proceedings of VLDB*, 1997.
- [9] M. Consens, D. Barbosa, A. Teisanu, and L. Mignet. Goals and benchmarks for autonomic configuration recommenders. In *Proceedings of SIGMOD*, 2005.
- [10] B. Dageville et al. Automatic SQL Tuning in Oracle 10g. In *Proceedings of VLDB*, 2004.
- [11] S. Papadomanolakis and A. Ailamaki. An integer linear programming approach to database design. In *Workshop on Self-Managing Database Systems*, 2007.
- [12] G. Valentin, M. Zuliani, D. Zilio, G. Lohman, and A. Skelley. DB2 advisor: An optimizer smart enough to recommend its own indexes. In *Proceedings of ICDE*, 2000.
- [13] C. Wu et al. The protein information resource: an integrated public resource of functional annotation of proteins. In *Nucleic Acids Research*, 2002.
- [14] D. Zilio et al. DB2 design advisor: Integrated automatic physical database design. In *Proceedings of VLDB*, 2004.
- [15] D. Zilio et al. Recommending materialized views and indexes with IBM DB2 design advisor. In *International Conference on Autonomic Computing*, 2004.

Nulls, Three-Valued Logic, and Ambiguity in SQL: Critiquing Date's Critique

Claude Rubinson*
Department of Sociology
University of Arizona
rubinson@u.arizona.edu

Abstract

Date's popular critique of SQL's three-valued logic [4, 3] purports to demonstrate that SQL queries can produce erroneous results when nulls are present in the database. I argue that this critique is flawed in that Date misinterprets the meaning of his example query. In fact, SQL returns the correct answer to the query posed; Date, however, believes that he is asking a different question. Although his critique is flawed, I agree with Date's general conclusion: SQL's use of nulls and three-valued logic introduces a startling amount of complexity into seemingly straightforward queries.

1 Introduction

A common critique of SQL is that the inclusion of nulls breaks the relational model. Date enumerates a number of reasons for this position. Most fundamentally, Date argues that—since SQL defines null not as a value but a flag indicating that the value of a particular attribute is missing—domains cannot properly include nulls since domains are, by definition, sets of values. Therefore, relations that include nulls are not, in fact, relations, undermining the very foundation of the relational model [3]. Date also make a more accessible argument in which he contends that the three-valued logic incurred by the use of nulls can generate

nonsensical results. In this essay, I critique this second argument and demonstrate that Date misapplies SQL's three-valued logic. Consequently, the critique is logically flawed and does not, in fact, indict SQL as Date supposes. Note, however, that my critique of Date is not a defense of nulls or SQL's three-valued logic; rather, it underscores just how confusing three-valued logic is. The introduction of nulls alters the meaning of seemingly straightforward queries and is likely responsible for numerous errors, errors which may frequently go unrecognized.

2 Date's Critique

Date's most prominent critique of nulls employs the simple SQL database illustrated in Figure 1. There are two tables. The Suppliers table (S) has two columns: the supplier number (the primary key) and the supplier's city. The Parts table (P) also has two columns: the part number (the primary key) and the part's city. In Figure 1, each table has only one record. Supplier S1 is located in London. We do not know in which city Part P1 is located.¹

¹Nulls often introduce confusion when it is unclear why information is missing from the database. Among the more common reasons for incomplete data entry are that the value of the attribute is (temporarily) unknown or that the attribute, itself, is not applicable to the represented entity. With regard to the present example, Date's discussion of the database described in Figure 1 makes it clear that the NULL marker in Table P indicates that the city associated with Part P1 is (temporarily) unknown. I proceed with this premise. In the conclusion, I return to this topic and discuss the additional complications

*I would like to thank Garrett Hoxie and Rick Snodgrass for their advice and support of this paper. I also wish to thank the *SIGMOD Record* reviewers for their helpful comments.

S	SNO*	CITY	P	PNO*	CITY
	S1	London		P1	NULL

Figure 1: SQL Database

Date [4, page 54] seeks to demonstrate that SQL’s three-valued logic produces erroneous results:

The fundamental point I want to make is that certain boolean expressions—and therefore certain queries—produces results that are correct according to three-valued logic but *not* correct in the real world.

To do so, he poses the following query: “Get SNO-PNO pairs where either the supplier and part cities are different or the part city isn’t Paris (or both)” [4, page 54] and writes the corresponding SQL implementation of the query:

```
SELECT S.SNO, P.PNO
FROM S, P
WHERE S.CITY <> P.CITY
OR P.CITY <> 'Paris'
```

Substituting in the data from the mock database, the expression (S.CITY <> P.CITY) OR (P.CITY <> 'Paris') becomes ('London' <> NULL) OR (NULL <> 'Paris') which, in accordance with the rules of SQL’s three-valued logic, evaluates to (NULL OR NULL) which, in turn, reduces to NULL. The query, therefore, returns no records.

Date [4, page 55] contends that this result reveals a flaw in SQL’s three-valued logic, arguing that:

But of course part P1 does have *some* corresponding city in the real world; in other words, “the null CITY” for part P1 does stand for some real value, say *xyz*. Obviously, either *xyz* is Paris or it isn’t.

Date then demonstrates that the WHERE clause will always evaluate to TRUE, regardless of where part P1 is located. In essence, there are three possibilities: city *xyz* is Paris, London, or some other city. If city *xyz*

that arise when the meaning of a null is ambiguous.

is Paris, the above expression becomes ('London' <> 'Paris') OR ('Paris' <> 'Paris'). This expression evaluates to (TRUE OR FALSE) which, in turn, evaluates to TRUE. If city *xyz* is London, the expression becomes ('London' <> 'London') OR ('London' <> 'Paris') which evaluates to FALSE OR TRUE which evaluates to TRUE. If city *xyz* is some other city, for example, New York, the expression becomes ('London' <> 'New York') OR ('New York' <> 'Paris') which evaluates to (TRUE OR TRUE) which, again, reduces to TRUE.

According to Date, if SQL correctly took account of the real world—specifically, that part P1 is associated with some city, despite that this fact is missing from the database—it should return the pair S1-P1. That SQL returns an empty set indicates a flaw in its logic: “In other words, the result that’s correct according to the logic (that is, 3VL) and the result that’s correct in the real world are different!” [4, page 55].

3 Critiquing the Critique

But Date is mistaken. The problem is not that SQL’s results disagree with reality but, rather, that Date poorly formulated his original inquiry. Recall Date’s original query: “Get SNO-PNO pairs where either the supplier and part cities are different or the part city isn’t Paris (or both).” The formulated SQL statement does not, in fact, correspond to this query; in fact, Date’s query cannot properly be translated into SQL because it assumes conventional, two-valued logic while SQL operates with three-valued logic.

In conventional logic, propositions are true or false. That is, part P1 is in Paris or it is not. In the three-valued logic employed by SQL, propositions are true, false, or unknown. By introducing the possibility of unknown propositions, it is no longer the case that part P1 is or is not in Paris. Rather: we know that part P1 is in Paris, we know that part P1 is not in Paris, or we don’t know where part P1 is.

The logic system within which one works demands that queries be formulated appropriately. Within a conventional, two-valued logic system, statements must be able to be classified as “true” or “false.”

Within a three-valued logic system, statements must also permit a classification of “unknown.” Date’s original query assumes two-valued logic. Consider the first clause of the query: “Get SNO-PNO pairs where . . . the supplier and part cities are different.” This query assumes that supplier and part cities are different or that they are the same. But within SQL’s three-valued logic system, supplier and part cities may be the same, they may be different, or we might not know if they are the same or different. The second clause of the query is similar: “Get SNO-PNO pairs where . . . the part city isn’t Paris.” Again, this query assumes that the part city is or is not Paris. Within SQL, however, the part city may be Paris, it may not be Paris, or we might not know what city it is. And, in fact, the null value in the database indicates that we do not know which city is associated with part P1.

Date argues that “in the real world” the city for part P1 either is or is not Paris. This is certainly true. But it is also true that “in the real world” we may not know what city is associated with part P1. These are two different propositions. The cities to which parts correspond is a set of facts that is distinct from whether we know which cities correspond to which parts. In SQL, queries always imply knowledge of the relationship in question and not simply the existence of said relationship. We can therefore reformulate Date’s original query as “Get SNO-PNO pairs where either we know that the supplier and part cities are different or we know that the part city isn’t Paris (or both).” The results of the SQL statement now make sense. An empty set is returned because—even though part P1 “does have *some* corresponding city in the real world”—we do not know to which city the part corresponds.

This understanding of the incongruity between two-valued and three-valued logic is made more clear by examining Date’s second example. Date [4, page 55] presents the following SQL statement

```
SELECT P.PNO
FROM   P
WHERE  P.CITY = P.CITY
```

and contends that “The real-world answer here is obviously the set of part numbers currently appearing

in P.” What is obvious is that Date thinks that the above SQL syntax is functionally equivalent to the statement “Get the PNO numbers for the parts that are associated with cities.” Because “in the real-world” all parts must be associated with a city, Date concludes that the query should return a set of part numbers. But Date is again misreading the query. Because SQL uses three-valued logic the statement expresses a distinctly different query: “Get the PNO numbers for the parts for which we know the associated city.” Again, SQL correctly returns an empty set because, according to table P, we do not know which city is associated with part P1.

Date [4, page 55] contends that his examples demonstrate that SQL is fundamentally broken:

To sum up: if you have any nulls in your database, you’re getting wrong answers to some of your queries. What’s more, you have no way of knowing, in general, just which queries you’re getting wrong answers to; *all* results become suspect. *You can never trust the answers you get from a database with nulls.* In my opinion, this state of affairs is a complete showstopper. (Emphasis in original.)

I have shown that Date has not demonstrated what he thinks he has. SQL returns the correct answer for the query posed but Date believes that he is asking a different question. This confusion is understandable. SQL’s three-valued logic is not intuitive. We are used to two-valued logic in which propositions are either true or false. But three-valued logic also permits unknown propositions. When working with SQL databases, it is imperative that we formulate our queries correctly; otherwise, we risk making mistakes similar to Date.

4 Discussion

SQL’s use of three-valued logic and its inclusion of the null marker requires that we formulate our database queries to reflect the possibility that the relationships between entities may be unknown. When we fail to

do so, we risk posing a different question than intended. We must keep in mind that SQL’s logic is non-intuitive. Rarely will the questions we put to a SQL database approach what we would ask in normal conversation. We cannot simply ask for the “SNO-PNO pairs where the supplier and part cities are different;” rather, we must ask for “SNO-PNO pairs where the supplier and part cities are known to be different.” More crucially, we must understand the difference between these two formulations.

The problem is only aggravated by the fact that information can be missing from a database for a variety of reasons. Date [1] identifies seven common causes of incomplete data entry: *value not applicable*, *value unknown*, *value does not exist*, *value undefined*, *value not valid*, *value not supplied*, and *value is the empty set*. If a value might be missing due to, for example, an inapplicable attribute, queries must be formulated and interpreted in consideration of this potential condition. When null markers are loaded with multiple meanings, the construction of associated queries rapidly becomes unmanageable: “Get SNO-PNO pairs where the part city attribute is applicable and either we know that the supplier and part cities are different or the part city isn’t Paris (or all three conditions apply).” To address this latter situation, some practitioners advocate the use of descriptive truth values [5, 7]. Constituting actual values rather than null markers, such solutions permit designers to construct databases that do not permit nulls and, consequently, may be queried using conventional, two-valued logic.

It may also be useful to note that any query that assumes three-valued logic may be decomposed into two correlated queries assuming two-valued logic.² Take, for example, the query “Get SNO-PNO pairs where we know that the part city isn’t Paris.” As discussed above, this query assumes three-valued logic because, for any given SNO-PNO pair, the part city may be in Paris, it may not be in Paris, or the part city may be unknown. This query may be decomposed into the compound query “Get SNO-PNO pairs where we know the part city and, from the resulting set, get SNO-PNO pairs where the part city is not Paris.” It

²I thank Charles Ragin for clarifying this principle for me.

is often helpful to perform this decomposition, particularly when constructing complex queries.

Ultimately, I agree with Date that three-valued logic is incompatible with database management systems. While I am not convinced that three-valued logic violates the relational model *per se*, I agree with McGoveran [6, page 355] that

many-valued logic means that database designers, developers, and users must all learn a whole new way of thinking. The practical costs of this approach are hard to assess; certainly they do violence to the goals we set out to satisfy with an RDBMS.

That Date, himself, misinterprets the meaning of his SQL syntax underscores the severity of the problem.

5 Conclusion

We develop databases in order to organize and make sense of information. The problem is that the world is complex. One manifestation of this complexity is that we sometimes lack complete information. I echo those who suggest that SQL practitioners avoid nulls as much as possible. By default, database designers should constrain columns as non-nullable. Operations that generate nulls such as outer joins should be avoided when possible, particularly as the basis of views and subqueries. Since, by definition, nulls indicate exceptional circumstances, nullable columns often indicate where the database design might be improved. The use of nulls in SQL is not the most fundamental concern raised by the database presented in Figure 1. Rather, it is: Where the heck is part P1? If part P1 is in transit to Paris, that information needs to be recorded in the database. So too if part P1 is lost. Notably, inclusion of such information elsewhere in the database increases both the value of the database as well as its integrity by permitting the problematic record to be dropped.

Proper design techniques, then, naturally minimize the number of nulls in the database. A database design is a model of a particular domain and it is only by thoroughly interrogating that domain—by

circumscribing its boundaries, delineating its constituent components, and identifying the relationships therein—that one can produce an accurate representation. Of course, the goal of a database design is not to represent a domain perfectly but only those aspects that are salient to the problem at hand. If part P1 is on a truck bound for Paris, its projected arrival time is probably relevant; that the truck driver just had a fight with his spouse, probably not. Nulls permit us to simplify our models by generalizing across anomalies that produce missing data and unknown relationships. But the cost of this simplified representation is three-valued logic and the associated increase in the complexity of our queries.

It is rare that one can guarantee the complete absence of nulls from a database. Even if database vendors were persuaded to deprecate nulls and three-valued logic, we would remain saddled with them for the foreseeable future. And since the presence of a single null value taints the entire database [6], one must generally assume three-valued logic. Consequently, the burden is on us to carefully review our queries to ensure that they mean what we intend.

References

- [1] C. J. Date. Not is not ‘not’! (notes on three-valued logic and related matters). In *Relational Database Writings, 1985–1989*. Addison Wesley, 1990.
- [2] C. J. Date. *Relational Database Writings, 1994–1997*. Addison Wesley Longman, 1998.
- [3] C. J. Date. *An Introduction to Database Systems*. Addison Wesley Longman, Reading, MA, seventh edition, 2000.
- [4] C.J. Date. *Database in Depth: Relational Theory for Practitioners*. O’Reilly, Sebastopol, CA, 2005.
- [5] G. H. Gessert. Handling missing data by using stored truth values. *SIGMOD Record*, 20(3):30–42, Summer 1991.
- [6] David McGoveran. Nothing from nothing (part 2 of 4) classical logic: Nothing compares 2 u. In *Relational Database Writings, 1994–1997* [2], chapter 6, pages 347–365.
- [7] David McGoveran. Nothing from nothing (part 4 of 4): It’s in the way that you use it. In *Relational Database Writings, 1994–1997* [2], chapter 8, pages 377–394.

A Data-oriented Survey of Context Models*

Cristiana Bolchini, Carlo A. Curino, Elisa Quintarelli, Fabio A. Schreiber, Letizia Tanca
Dip. di Elettronica e Informazione – Politecnico di Milano
P.zza Leonardo da Vinci, 32 – 20133 Milano (Italy)
{bolchini,curino,quintare,schreiber,tanca}@elet.polimi.it

ABSTRACT

Context-aware systems are pervading everyday life, therefore context modeling is becoming a relevant issue and an expanding research field. This survey has the goal to provide a comprehensive evaluation framework, allowing application designers to compare context models with respect to a given target application; in particular we stress the analysis of those features which are relevant for the problem of data tailoring. The contribution of this paper is twofold: a general analysis framework for context models and an up-to-date comparison of the most interesting, data-oriented approaches available in the literature.

1. INTRODUCTION

Many interpretations of the notion of context have emerged in various fields of research like psychology, philosophy [13], or computer science [31]. Context has often a significant impact on the way humans (or machines) act and on how they interpret things; furthermore, a change in context causes a transformation in the experience that is going to be lived. The word itself, derived from the Latin *con* (with or together) and *texere* (to weave), describes a context not just as a profile, but as *an active process dealing with the way humans weave their experience within their whole environment, to give it meaning*.

While the computer science community has initially perceived the context as a matter of user location, as Dey and Abowd discuss in [2], in the last few years this notion has been considered not simply as a state, but part of a process in which users are involved [18]; thus, sophisticated and general context models have been proposed, to support context-aware applications which use them to (a) adapt interfaces [20], (b) tailor the set of application-relevant data [8], (c) increase the precision of information retrieval [43], (d) discover services [40], (e) make the user interaction implicit [37], or (f) build smart environments [21].

Accordingly, consider the example of automated support for a natural history museums visitors, who may be endowed with a portable device which reacts to a change of context by (a) adapting the user interface to the different abilities of the visitor – from low-sighted people to very young children –; (b) providing different information contents based on the different interests/profiles of the visitor (geology, paleontology, ... scholar, journalist, ...), and on the room s/he is currently in; (c) learning, from the previous choices per-

formed by the visitor, what information s/he is going to be interested in next; (d) providing the visitor with appropriate services – to purchase the ticket for a temporary exhibition, or to reserve a seat for the next in-door show on the life of dinosaurs –; (e) deriving location information from sensors which monitor the user environment; (f) provide active features within the various areas of the museum, which alert visitors with hints and stimuli on what is going on in each particular ambient.

Artificial Intelligence developed, since the late 80s, a notion of context [23, 25, 34, 35, 42] that differs from the one considered in this paper. The AI goal was extending the existing reasoning techniques to enable contextual reasoning. The most mature approaches are Propositional Logic of Context and MultiContext System/Local Models Semantics. While the first introduces the context as a “first class citizen” of a logic theory, the second perceives context as “a partial and approximate theory of the world from some individual’s perspective”. Both succeed in modeling context to enable reasoning and provide extremely expressive mechanisms to exploit context in formal theories, as proved by their recent application to the Semantic Web [11, 26]. However the need for a simple, explicit, unified model of context able to gather in a single representation several individual contexts, requires a rather different approach, whose features are presented and analyzed in Section 2.

In the general high-level architecture of a context-aware system, context design is carried out according to the application domain, by modeling the elements that affect the knowledge/services/actions that have to be made available to the user at run-time, when a context becomes active.

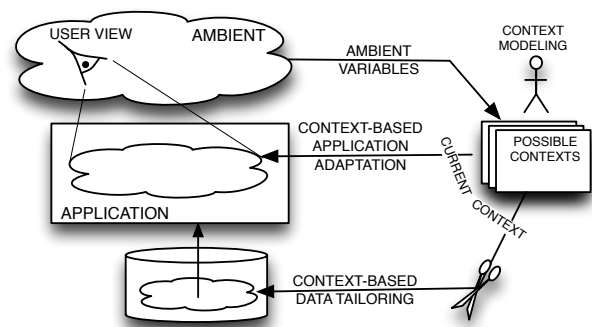


Figure 1: A context-aware system architecture

*This research is partially supported by the Italian MIUR projects: ART-DECO (FIRB), and ESTEEM (PRIN).

The context information acts as the command source for input- and output- related switches, which enact alternate behaviors providing different information while all the rest remains unchanged. While in a traditional system context data are not treated as special information and the system implicitly includes all different behaviors without being aware of the multiple facets of the application ambient, in a context-aware system, context data are used to customize the way inputs are processed (Figure 1).

In Information Management, context-aware systems are mainly devoted to determining *what portion of the entire information is relevant with respect to the ambient conditions*.

Given this scenario, *context-based data tailoring* [10] can be defined as the activity of *defining data views*, based on a) the identification of the various contexts the application user is going to experience in the envisaged scenario, b) the design of a set of data views for each of the identified contexts. The aim is to provide support to the designer of data management applications, be them related to a huge (e.g., in data warehousing) or to a very small amount of data (e.g., in portable, lightweight data management systems), in determining and creating the various views to be used in the different contexts, by following a systematic approach. Indeed, nowadays the amount of available data and data sources requires not only to integrate them (still a hard problem), but also to filter (tailor) the relevant portion of data in order to: 1) provide the user with the appropriately tailored set of data, 2) match devices' physical constraints, 3) operate on a manageable amount of data (for improving query processing efficiency), and 4) provide the user with time- and location-relevant data (mobile applications).

We select the data tailoring issue as our target application because we consider it as an enabling component for the forthcoming Information Systems, such as mobile, data-intensive systems, P2P systems and in general the Semantic Web. In particular, in the last research area, huge ontologies (several millions of concepts and relations) are starting to appear, e.g., UMLS [32], while common operations such as query answering, reasoning and consistency check may be exponential in the size of the input ontology [4]¹. Transparent context-aware sub-ontology extraction, exploiting techniques similar to the one proposed in [6], can improve the performance of ontology manipulation, while preserving the user perception of operating on the complete ontology.

Interesting surveys on context-aware systems and models have already been presented in [5, 14, 39, 45]; we contribute with a review on recent evolutions and new systems for context modeling, and challenge each model with respect to the problem of context-aware data tailoring.

We believe that no “silver-bullet” in context modeling has been proposed so far, and that a deep understanding of the context problem itself is essential to choose or design the right model; for this purpose we introduce in Section 2 a framework useful for analyzing context models and to select the most suitable one for a given application. In fact, the lack of a uniform approach for modeling context-

¹Many of these tasks may have lower complexity for logics which are not very expressive, like for instance \mathcal{FL}^- [4]. However, most of the interesting ontologies found on the Web are at least as expressive as OWL-Lite (*SHIF*), where, for instance, just concept satisfiability is EXPTIME-complete [29].

related information makes it difficult to deeply understand the requirements that have to be considered when proposing/adopting a context model on the basis of its focus. Therefore, the central issue of this paper is to survey the current literature on the context modeling problem and to systematically highlight advantages and limitations of the different proposals and perspectives.

The rest of the paper is organized as follows, Section 2 describes the analysis framework, Section 3 applies it to some of the most relevant context models found in the literature, Section 4 draws some conclusions.

2. THE ANALYSIS FRAMEWORK

Many approaches defining the notion of context have been proposed and several adaptive applications have been designed and implemented, by introducing the notions of user profile and context [1, 3, 9, 12, 13, 19, 33, 36, 47]. Although interesting comparisons of context models already exist [5, 30, 39, 45], we felt the need to establish a framework to systematically evaluate them, by defining a set of relevant, objective and rather general categories. The analysis framework we propose is intended for designers that are about to develop context-aware applications and need to decide which context model is best suited for their goals. This framework, used to analyze and compare the available context models, is built on a rich set of features which characterize the models from various perspectives. These features have been derived from the analyzed systems, by selecting the most peculiar and common ones. The first step of the analysis is the identification of the key issues for the application being developed; in this phase the designer should define which features are more relevant for his/her target application, or whether new features should be added to address specific application requirements. Here we assume *data tailoring* as our target application and, with respect to it, we show the most relevant features among the presented ones. The second step is the classification of the existing context models with respect to each feature. The result is a structured view of the state of the art, which enables the designer to consciously compare the various models, focusing the attention on the key issues isolated in step one. As a result, the best model is selected or, in case no satisfactory models are available for the target application, the designer might consciously engage in the proposal of a new, more appropriate context model. The features we isolated and classified are now briefly discussed:

Modeled aspects: The set of context dimensions managed by the model.

- *Space*: does the considered context model deal with location-related aspects?
- *Time*: does the considered context model allow the representation of temporal aspects?
- *Absolute/relative space and time*: are the space and time parameters (if any) represented absolutely (e.g., GMT time reference and GPS coordinates) or relatively (e.g., “near something”, “last month”, “after that”)?
- *Context history*: is the history of previous contexts part of (relevant for) the context itself, i.e., the current context state depends on previous ones, or is the

context a pure snapshot of the user's current environment?

- *Subject*: who or what is the subject of the described context? This feature refers to the point of view used to describe the context itself; some models describe the context as it is perceived by the user, while others assume the application point of view, considering, as a consequence, the user itself as part of the context;
- *User profile*: is the user profile (in terms of preferences and personal features) represented in the context model? And if so, how is it represented (i.e., does the system describe the user's characteristics one by one, or does it provide a role-based model of user classes)?

Representation features: General characteristics of the model itself.

- *Type of formalism*: class of the conceptual tool used to capture the context (key-value-, mark-up scheme-, logic-, graph-, ontology-based). Different classes provide different features (e.g., high or low intuitiveness, possibility to be automatically processed, reasoning support, formal semantics) and are more or less adequate for certain applications;
- *Level of formality*: the existence of a formal definition and whether the formalization well expresses the intuition;
- *Flexibility*: the model's ability to easily adapt to different contexts: a model can be "application-domain bounded" if it is substantially focused on a single application or on a specific domain, or "fully general" if it can naturally deal with different domains or applications (i.e., is it possible to capture any kind of context with this model and how easy is it?);
- *Variable Context Granularity*: the ability of the model to represent the characteristics of the context at different levels of detail.
- *Valid Context Constraints*: the possibility to reduce the number of admissible contexts by imposing semantic constraints that the contexts must satisfy for a given target application.

Context management and usage: The way the context is built, managed and exploited.

- *Context construction*: highlights if the context description is built centrally or via a distributed effort; this indicates whether a central, typically design-time, description of the possible contexts is provided, or if a set of partners reaches an agreement about the description of the current context at run-time;
- *Context reasoning*: indicates whether the context model enables reasoning on context data to infer properties or more abstract context information (e.g., deduce user activity combining sensor readings);
- *Context information quality monitoring*: indicates whether the system explicitly considers and manages the quality of the retrieved context information, for instance, when the context data are perceived by sensors;

- *Ambiguity and incompleteness management*: in case the system perceives ambiguous, incoherent or incomplete context information, indicates if the system can "interpolate" and "mediate" somehow the context information and construct a reasonable "current context";
- *Automatic Learning Features*: highlights whether the system, by observing the user behavior, individual experiences of past interactions with others, or the environment, can derive knowledge about the context; e.g., by studying the user's browsing habits, the system learns user preferences;
- *Multi-Context Modeling*: the possibility to represent in a single instance of the model all the possible contexts of the target application, as opposite to a model where each instance represents a context.

This characterization covers the focus of the model, its representation and the way context data are used; the result is a rich set of features, emphasizing that context modeling is a varied and complex problem. Depending on the specific purpose it is designed for, each model may "include" several of the listed features; we envision five classes of use, which share general sets of features, and more important, the same target field of application. These classes can be considered as a coarse-grained categorization of the context models, or as a decomposition of the context problem itself (in boldface the key features of each class).

- Context as a matter of channel-device-presentation.* Systems of this class are characterized by: **variable context granularity**, the **application as subject** of the model, limited or absent management of location and time dimensions, **feature-based user profiling**, low level of formality, limited flexibility (often considering only specific applications), and a **centrally defined context**. While automatic learning features can be available, context quality monitoring, ambiguity management and context reasoning are in general not supported.
- Context as a matter of location and environment.* Models of this class in general provide: precise **time and space management**, high degree of flexibility and **centralized context definition**. Context reasoning may be provided, offering a powerful abstraction mechanism. **Information quality management** and **disambiguation** may be available, in particular when the context information is acquired by sensors. Automatic learning is rarely exploited.
- Context as a matter of user activity.* The focus of this class of models is on "what the user is doing," consequently **context history** and **reasoning** are important issues. Time and space are considered relevant as far as they provide information about the user current activity². While the level of formality may vary, the **context definition** is in general **centralized** and the **user** is the **subject of the model**. When available, the **automatic learning** is used to guess user activity from sensor readings.

²See [37] for an example

- D. *Context as a matter of agreement and sharing (among groups of peers)*. Approaches of this group focus on the problem of reaching an agreement about a context shared among peers; clearly the **context definition is distributed**; **context reasoning**, **context quality monitoring** and **ambiguity and incompleteness management**, are key issues. Sophisticated location, time and user profiling features are uncommon in models of this class. The **level of formality** is rather **high**, due to the need of information sharing.
- E. *Context as a matter of selecting relevant data, functionalities and services (data or functionality tailoring)*. The models of this group focus on how the context determines which data, application functionalities and services are relevant. Context definition is typically centralized, context history and reasoning are often not provided; **time**, **space** and **user profile** are in general highly developed and well formalized. The flexibility is usually high while automatic learning features, ambiguity management and information quality are not key issues and are often not available. The key features of this group are: **the application as subject**, the possibility to express both **variable context granularity**, **valid context constraints**, and **multi-context models**.

These classes and the identified relevant features constitute the analysis framework we propose, used in the next section to review some of the most interesting approaches to the context modeling problem.

3. THE CONTEXT MODELS

Table 1 reports the results of the application of the analysis framework to a set of systems examined with the data tailoring application scenario in mind³. A very short description of each system follows, highlighting relevant characteristics and the context modeling subproblems they are targeting.

- *ACTIVITY*³: in [30] the authors provide an interesting analysis of the existing approaches to context modeling, pointing out how different solutions overlap without providing the context modeling universal solution. The authors also describe a novel approach based on Activity Theory, which allows the description of key aspects influencing human activity. In fact, in [30] the notion of context is intended as the set of elements which have some influence on users' intentions while performing an activity. The model is strongly focused on the categories of *user*, *community* and the *rules* needed to relate a user to his/her community; each category can be represented by a tree-based structure, where lower levels of the tree represent more detailed information about the context category that can be used for reasoning about upper levels. To the best of our knowledge, a formal description of the context model has not been provided and its usage is not described; the model seems to be at a very early stage of development, and too holistic to be effective in practice. The ultimate goal is the context problem as a whole, fitting into all our categories.

³Each model appears in all the applicable categories, possibly with more values per category.

- *CASS*: it is a centralized server-based context management framework, meant for small portable devices, offering a high-level abstraction on context sensed by appropriate distributed sensors [22]. It manages both time and space, taking into account the context history, and provides context reasoning; it does not contain user profiling capabilities. The context is seen as a matter of location and environment, thus the system can be classified into the B category.
- *CoBrA*: The context is represented as a Context Knowledge Base [16] for the specific application of event/meeting management. On top of this knowledge base temporal, spatial and event-meeting reasoners (based on contextual rules) operate to deduce more abstract contextual information. The presence of a Context Broker makes this approach perfectly suited for context sharing and context reasoning, while its application is difficult when multiple multidimensional contexts need to be modeled. To apply CoBrA [15] to the information tailoring problem we must enrich the ontologies forming the Context Knowledge Base to extend the domain of applicability from the "meeting" domain to other application-specific domains and to define a set of contextual rules describing how various components of the context should be combined. Such rules, although specifiable as CoBrA context reasoning rules, will express how to combine context characteristics instead of supporting contextual inference, therefore forcing the original model to suit this specific goal. CoBrA belongs to the D group.
- *CoDaMoS and SOCAM*: The CoDaMoS [17, 38] and the SOCAM projects [24] propose extremely general ontology-based context models. Sets of extensible ontologies are exploited to express contextual information about user, environment and platform in both systems. CoDaMoS is the two-layered context model used in PACE [27], a middleware for context aware systems, which describes contexts both in term of fine grained facts and higher level situations which describe logical conditions; CoDaMos adds also support for service description. The richness and flexibility of such models is not complemented by a proper constraining mechanism; the two models do not offer explicit ways to limit the number of expressible contexts (i.e., Valid Context Constraints), this results in a severe limitation when the context model is applied to the data tailoring problem. Moreover a single point in the multidimensional context space is not represented in a concise way, but as a graph of concept instances, making the task of relating the set of relevant data to the specific context difficult. The possibility to express contexts at different granularity levels and to define them compositionally (e.g., as a combination of more detailed ones) is also difficult to achieve (i.e., Variable Context Granularity). Both participate to the four above-mentioned classes, being more focused in A and B.
- *COMANTO*: [41,46] the authors propose a hybrid context modeling approach to handle context objects and context knowledge. For the first purpose, a location-based context model is formalized for considering both fixed (e.g., regions, streets, etc.) and mobile location

System	Space	Time	Space/Time coordinates (R elative or A bsolute)	Context history	Subject (U ser or A pplication)	User profile (R ole or F eatures based)	Variable context granularity	Valid context constraints	Type of formalism: Key-value-based	Type of formalism: Mark-up based	Type of formalism: Logic-based	Type of formalism: Graph-based	Type of formalism: Ontology-based	Formality level (H igh or L ow)	Flexibility	Context construction (D istributed or C entralized)	Context reasoning	Context quality monitoring	Ambiguity/Incompleteness mgmt.	Automatic learning features	Multi-context model
ACTIVITY	+		A	+	U	F	+					+		L	+	C	+				+
CASS	+	+		+	U	F					+			L		D	+				
CoBrA	+	+	A		A	F								H	+	D	+		+		
CoDaMoS	+	+	R/A		A	F								H	+	D	+	+		+	
COMANTO	+	+	R/A		A	F								H	+	D	+				
Context-ADDICT	+	+	R/A		A	R	+	+		+		+	+	H	+	C	+				
Conceptual-CM	+	+	R	+	A	R						+		L	+	C	+			+	+
CSCP					A	F				+				L		C			+		+
EXPDOC		+	R	+	U	F					+	+		H		C				+	+
FAWIS				+	U	F	+		+					H		C			+	+	+
Graphical-CM	+	+	R		A	F						+		H	+	C	+	+		+	+
HIPS/HyperAudio	+	+	A	+	U	F			+					L		C	+		+	+	+
MAIS	+	+	A		A	F	+							H		C			+	+	+
SCOPEs					A	F					+			H		D	+		+		+
SOCAM	+	+	R/A		A	F							+	H	+	D	+	+			
U-Learn	+	+	A		U	F	+						+	H	+	D					+

Table 1: Context model features and systems exposing them.

data (e.g., people, vehicles). For the second purpose the general COMANTO ontology is proposed as a public context semantic vocabulary supporting efficient reasoning on contextual concepts (such as users, activities, tools, etc.) and their associations. The ontology is used to collect a structured semantic representation about generic context information and is not domain-, or application-oriented. The middleware infrastructure to acquire, store, and manage context information of the COMANTO ontology is described in [46]. As for the other context models based on ontologies, COMANTO provides a general purpose and very expressive formal model, although lacking the possibility to discard useless contexts. This model fits into categories B and C.

- *ConceptualCM*³ [18]: ConceptualCM is a conceptual framework intended to consider the context notion not simply as a state, but as part of a process. The possible contexts for a scenario are an information space modeled as a directed state graph, where each node represents a context and edges denote the conditions for changing context. Each context is defined by a set

of entities, a set of roles that entities must satisfy, a set of relations between entities, and a set of situations. A runtime infrastructure is a middleware that instantiates entities, roles and relations for the current state of a context, with different levels of abstraction, by allowing the collection of all the information required to identify current context values and predict changes in the situation or in the actual context. In [18] the authors describe some principles to be considered when implementing context-aware applications; the context model informally presented can be classified into the E category.

- *Context-ADDICT*: In [7] the authors propose the Context Dimension Tree, a tree-based structure introduced in the research field of context aware applications with the specific goal of being adopted in the data tailoring task. The model globally represents the space of considered contexts; in particular, the root node of the tree specifies the entire data space of possible contexts, and the first-level nodes (called dimensions) represent the orthogonal perspectives to be considered in order to tailor data. The hierarchical structure of the Con-

text Dimension Tree can represent contexts with different levels of detail, and the portion of data to be considered for a specific context can be determined in a compositional way, by using the data relevant for each dimension value composing the current context. The model includes constraints and relationships among dimension values to remove meaningless combinations of elements. Being focused on a specific class of applications, the Context-ADDICT approach lacks the features not relevant for the data tailoring problem such as Context History, Context Quality Monitoring, Context Reasoning and Ambiguity and Incompleteness Management. Some of these limitations may be removed in the future while other are inherent to the chosen approach. This model is classified as pertaining to the E category.

- *CSCP*: in [12] the authors present a Mobility Portal: a web portal providing an adaptive web interface, reacting to user channel, device and user profile. The focus is clearly on channel-device-presentation issues, thus the contribution is limited to a well defined set of applications, based on web interfaces. The context model represents profile sessions and is based on RDF; it does not impose any fixed hierarchical structure for the context notion, thus inherits the full flexibility and expressive power of RDF. The instantiation of the model allows one to represent a single structured session profile (i.e., a point in the space of possible contexts) with information about the device, the network, and the user of the considered session. The best classification of this system is in group A.
- *EXPDOC*⁴: it is an interesting approach based on semantic networks [44]. The goal is to support experiential systems (in particular experiential documents), in order to provide an enriched learning environment where additional, related, but not required information is made available to the users, the authors talk about “serendipitous” activities, as the set of knowledge improving activities performed by the users when accessing this kind of information. The goal of this approach is opposite to the one of data tailoring: while EXPDOC uses the context to increase the amount of information provided to the user, data tailoring exploits the context to discard useless information. As a consequence, this semantic-network-based approach suffers from the same limitations we discussed for the ontological models, in particular there are no ways to limit the contexts expressed by this model (i.e., Valid Context Constraints). Moreover, the automatic Wordnet-based mechanisms, exploited to generalize the user context in order to match the document context, focus only on the user preferences and profiles, resulting in limited flexibility (e.g., the location is not taken into account). For these reasons, the application of such model to the data tailoring problem is not advisable. This system belongs to the C class.
- *FAWIS*: the methodology of [19, 20] is focused on the adaptation of Web-based Information Systems via the

⁴This is not the original name, it has been introduced to easily refer to this system in Table 1.

transformation of the presentation and navigation, although the context model is flexible enough to be applied to different scenarios. A specific context is specified by a set of profiles, each describing an autonomous aspect of the context itself (e.g., the *user*, the *location*, the *device*, etc.). A profile is characterized by a set of simple or complex attributes, and each instantiation of a profile has a fixed set of attributes, assuming also the presence of null values. Profiles can be combined to represent a context at different levels of detail; however, the model does not allow the expression of constraints between sets of attributes or set of profiles to avoid meaningless combinations. The system mainly considers the user-profiling issues of the context modeling problem, while leaving all the other aspects not formally described, thus, in our classification it falls into group A and partially into groups E and C.

- *GraphicalCM*³ [28]: the authors formalize a context model for pervasive computing applications, by concentrating also on some aspects not well formalized in the literature for this specific field, that are information quality and temporal aspects of contexts. The context model has a graphical notation: the possible contexts for a target application are rendered by a directed graph composed by a set of entities, describing objects, and their attributes, representing the entity properties. Different kinds of associations connect an entity to its attributes or to other entities. GraphicalCM supports quality by annotating associations with a number of quality parameters, which capture the related dimensions of quality considered relevant for each association. Each quality parameter is described by one or more quality metrics. The model is theoretically described in [28], and the authors introduce some possible extensions for their proposal, which is general, but at the moment can be classified into the C and E categories.
- *HIPS/HyperAudio*: the authors of [37] focus their attention on the spatio-temporal issues of the context, concentrating on determining the user’s current activity from information about his/her spatio-temporal coordinates and a simple user profile. They consider the context as a matter of user activity, and the target of their key-value-based model is supporting an automatic context-aware museum guide. Although the approach is rather effective in this specific application, it has limited flexibility. The exploitation of this model in more general applications (e.g., data tailoring) is definitely hard, and major extensions are required to capture articulated contexts. This approach may be considered as belonging to categories B and C.
- *MAIS* [9, 33]: it is a Multi Channel Adaptive Information System meant to build a flexible environment to adapt the interaction and provide information and services according to ever-changing requirements, execution contexts, and user needs. The notion of MAIS context has the objective of configuring the software on board of the device based on: a) the user needs, in terms of presentation, and b) the device characteristics, in terms of available channel. It clearly considers the context as a matter of channel-device-presentation,

thus it belongs to group A.

- *SCOPES*: the context model presented in [36] is based on the concept of *mutual beliefs*. In a P2P collaborative environment assertions are exchanged among peers to create mappings among source schemata. These sets of mapping represent the notion of evolving context proposed by the authors. The goal of the system is to enable P2P data interoperability, via the definition of the above described context. Although the system shares some of the goals typical of data tailoring systems, the presented context model cannot be applied in the data tailoring task: it is not possible with this mutual-beliefs-based approach to define a context model independent of the data sources, in fact it does not include constructs to represent elements like location or user profile. The model falls into class D.
- *U-Learn*³ [48]: this ontology-based context-model is focused on the support of learning. The learner and the learning content are described by two ontologies (learner ontology and content metadata) and a rule-based system provides a content-to-learner matching mechanism. The content can be both a service or a set of data. The proposal is interesting with respect to the data tailoring problem; the data can be enriched by adding content metadata, the user's context described by the learner ontology and the matching can be used to select the relevant data depending on the context. Yet, the system seems at an early stage of development, and the formalization not complete: the learner and learning-content ontologies seem very general and not clearly specified, while the matching rules are not described in the available papers. The authors claim to support sensor integration without providing enough details to actually evaluate the contribution. This model can be classified into the E category.

4. CONCLUSIONS

Although a lot of work has been done, the representation and management of context can hardly be considered as an assessed issue. Due to the complexity of the “context modeling problem” as a whole and to the multitude of different applications, at the end of this comparison we advocate those models that, although being fully general, have a well defined focus, and try to support only a specific context subproblem. Indeed, we feel that the systems whose aim is to be completely general and to support the context modeling problem as a whole for any possible application, often fail to be effective. In fact, the practical applicability and usability, although not discussed because rather subjective, are important parameters, and are often inversely proportional to the generality of the model: the more expressive and powerful, the less practical and usable.

Different context subproblems and applications have almost incompatible requirements, and common solutions are still not available; as a consequence, the context model should be chosen depending on the target application. The analysis framework we have proposed can, in this sense, be used by an application designer either to choose among the available models or to define the requirements of a new context model.

5. REFERENCES

- [1] K. Aberer and et al. Emergent semantics: Principles and issues. In *Invited paper at 9th Int. Conf. on Database Systems for Advanced Applications, LNCS 2973*, pages 25–38, 2004.
- [2] G. D. Abowd, A. K. Dey, P. J. Brown, N. Davies, M. Smith, and P. Steggles. Towards a better understanding of context and context-awareness. In *Proc. 1st Int. Symp. on Handheld and Ubiquitous Computing (HUC'99)*, pages 304–307. Springer-Verlag, 1999.
- [3] R. Agrawal and E. L. Wimmers. A framework for expressing and combining preferences. In *Proc. ACM SIGMOD Int. Conf. on Management of Data 2000*, pages 297–306. ACM, 2000.
- [4] F. Baader, D. Calvanese, D. McGuinness, D. Nardi, and P. Patel-Schneider, editors. *The Description Logic Handbook: Theory, Implementation, and Applications*. Cambridge University Press, 2003.
- [5] M. Baldauf, S. Dustdar, and F. Rosenberg. A survey on context-aware systems. *Int. Journal of Ad Hoc and Ubiquitous Computing*, 2(4):263–277, June 2007.
- [6] M. Bhatt, A. Flahive, C. Wouters, W. Rahayu, and D. Taniar. Move: A distributed framework for materialized ontology view extraction. *Algorithmica*, 45(Issue 3):457–481, 2006.
- [7] C. Bolchini, C. Curino, E. Quintarelli, F. A. Schreiber, and L. Tanca. Context-ADDICT. Technical Report 2006.044, Dip. Elettronica e Informazione, Politecnico di Milano, 2006.
- [8] C. Bolchini, C. Curino, F. A. Schreiber, and L. Tanca. Context integration for mobile data tailoring. In *Proc. 7th IEEE/ACM Int. Conf. on Mobile Data Management*, page 5, 2006.
- [9] C. Bolchini, F. A. Schreiber, and L. Tanca. Data management. In B. Pernici, editor, *Mobile Information Systems - Infrastructure and Design for Adaptivity and Flexibility*, chapter 6, pages 155–176. Springer, 2006.
- [10] C. Bolchini, F. A. Schreiber, and L. Tanca. A methodology for very small database design. *Information Systems*, 32(1):61–82, March 2007.
- [11] P. Bouquet, F. Giunchiglia, F. Harmelen, L. Serafini, and H. Stuckenschmidt. C-OWL: Contextualizing Ontologies. In *2nd Intl Semantic Web Conference, LNCS 2870*, pages 164–179, 2003.
- [12] S. Buchholz, T. Hamann, and G. Hübsch. Comprehensive structured context profiles (CSCP): Design and experiences. In *Proc. 2nd IEEE Conf. on Pervasive Computing and Communications Workshops*, pages 43–47, 2004.
- [13] M. Chalmers. A historical view of context. *Computer Supported Cooperative Work*, 13(3):223–247, 2004.
- [14] G. Chen and D. Kotz. A survey of context-aware mobile computing research. Technical report, Dept. of Computer Science, Dartmouth College, Hanover, NH, USA, 2000.
- [15] H. Chen, T. Finin, and A. Joshi. An Intelligent Broker for Context-Aware Systems. *Adjunct Proc. of Ubicomp 2003*, pages 183–184, October 2003.
- [16] H. Chen, F. Perich, T. Finin, and A. Joshi. SOUPA: Standard Ontology for Ubiquitous and Pervasive Applications. In *Int. Conf. on Mobile and Ubiquitous*

- Systems: Networking and Services*, August 2004.
- [17] CoDaMoS development team. The codamos project, 2003.
- [18] J. Coutaz, J. L. Crowley, S. Dobson, and D. Garlan. CONTEXT is KEY. *Communications of the ACM*, 48(3):49–53, 2005.
- [19] R. De Virgilio and R. Torlone. A general methodology for context-aware data access. In *Proc. ACM Int. Workshop on Data engineering for wireless and mobile access*, pages 9–15, 2005.
- [20] R. De Virgilio, R. Torlone, and G.-J. Houben. A rule-based approach to content delivery adaptation in web information systems. In *Proc. 7th IEEE/ACM Int. Conf. on Mobile Data Management*, page 21, 2006.
- [21] A. K. Dey, T. Sohn, S. Streng, and J. Kodama. iCAP: Interactive prototyping of context-aware applications. In *Proc. 4th Int. Conf. on Pervasive Computing*, pages 254–271, 2006.
- [22] P. Fahy and S. Clarke. CASS - middleware for mobile context-aware applications. In *Proc. Mobisys 2004 Workshop on Context Awareness*, 2004.
- [23] F. Giunchiglia. Contextual reasoning. *Epistemologia*, 16, 1993.
- [24] T. Gu, H. K. Pung, and D. Q. Zhang. A service-oriented middleware for building context-aware services. *Journal of Network and Computer Applications*, 28(1):1–18, 2005.
- [25] R. V. Guha and J. McCarthy. Varieties of contexts. In *Proc. 4th. International and Interdisciplinary Conference, CONTEXT, LNAI 2680*, pages 164–177, 2003.
- [26] R. V. Guha, R. McCool, and R. Fikes. Contexts for the semantic web. In *International Semantic Web Conference*, pages 32–46, 2004.
- [27] K. Henriksen, J. Indulska, T. McFadden, and S. Balasubramaniam. Middleware for distributed context-aware systems. In *Proc. of On the Move to Meaningful Internet Systems, LNCS 3760*, pages 846–863, 2005.
- [28] K. Henriksen, J. Indulska, and A. Rakotonirainy. Modeling context information in pervasive computing systems. In *Proc. 1st Intl Conf. on Pervasive Computing, LNCS 2414*, pages 167–180, 2002.
- [29] I. Horrocks, U. Sattler, and S. Tobies. Practical reasoning for very expressive description logics. *J. of the Interest Group in Pure and Applied Logic*, 8(3):239–264, 2000.
- [30] M. Kaenampornpan and E. O’Neill. An intergrated context model: Bringing activity to context. In *Proc. Workshop on Advanced Context Modelling, Reasoning and Management*, 2004.
- [31] D. Lenat. The dimensions of context-space. *CYCORP*, 1998.
- [32] D. A. Lindberg, B. L. Humphreys, and A. T. McCray. Unified medical language system project: A distributed experiment in improving access to biomedical information. *Methods in Inf. Med.*, 32(4):281–291, 1993.
- [33] MAIS. Mais: Multi channel adaptive information system.
- [34] J. McCarthy. Notes on formalizing context. In *IJCAI*, pages 555–562, 1993.
- [35] J. McCarthy and S. Buvač. Formalizing context (expanded notes). In *Computing Natural Language*, volume 81 of *CSLI Lecture Notes*, pages 13–50. 1998.
- [36] A. M. Ouksel. In-context peer-to-peer information filtering on the web. *SIGMOD Record*, 32(3):65–70, 2003.
- [37] D. Petrelli, E. Not, C. Strapparava, O. Stock, and M. Zancanaro. Modeling context is like taking pictures. In *Proc. of the Workshop "The What, Who, Where, When, Why and How of Context-Awareness" in CHI2000*, 2000.
- [38] D. Preuveneers, J. V. den Bergh, D. Wagelaar, A. Georges, P. Rigole, T. Clerckx, Y. Berbers, K. Coninx, V. Jonckers, and K. D. Bosschere. Towards an extensible context ontology for ambient intelligence. In *Proc. 2nd European Symp. Ambient Intelligence, LNCS 3295*, pages 148–159, 2004.
- [39] D. Raptis, N. Tselios, and N. Avouris. Context-based design of mobile applications for museums: a survey of existing practices. In *Proc. 7th Int. Conf. on Human computer interaction with mobile devices & services*, pages 153–160, 2005.
- [40] P.-G. Raverdy, O. Riva, A. de La Chapelle, R. Chibout, and V. Issarny. Efficient context-aware service discovery in multi-protocol pervasive environments. In *7th IEEE/ACM Int. Conf. on Mobile Data Management*, page 3, 2006.
- [41] I. Roussaki, M. Strimpakou, N. Kalatzis, M. Anagnostou, and C. Pils. Hybrid context modeling: A location-based scheme using ontologies. In *4th IEEE Conf. on Pervasive Computing and Communications Workshops*, pages 2–7, 2006.
- [42] L. Serafini and P. Bouquet. Comparing formal theories of context in AI. *Artificial Intelligence*, 155(1-2):41–67, 2004.
- [43] X. Shen, B. Tan, and C. Zhai. Context-sensitive information retrieval using implicit feedback. In *Proc. 28th Int. ACM SIGIR Conf. on Research and development in information retrieval*, pages 43–50, 2005.
- [44] H. Sridharan, H. Sundaram, and T. Rikakis. Computational models for experiences in the arts, and multimedia. In *Proc. ACM Workshop on Experiential Telepresence*, pages 31–44, 2003.
- [45] T. Strang and C. Linnhoff-Popien. A context modeling survey. In *1st Int. Workshop on Advanced Context Modelling, Reasoning and Management*, 2004.
- [46] M. Strimpakou, I. Roussaki, and M. E. Anagnostou. A context ontology for pervasive service provision. In *20th Int. Conf. on Advanced Information Networking and Applications*, pages 775–779, 2006.
- [47] R. Torlone and P. Ciaccia. Management of user preferences in data intensive applications. In *Proc. 11th Italian Symp. on Advanced Database Systems*, pages 257–268, 2003.
- [48] S. J. H. Yang, A. Huang, R. Chen, S.-S. Tseng, and Y.-S. Shen. Context model and context acquisition for ubiquitous content access in ulearning environments. In *IEEE Int. Conf. Sensor Networks, Ubiquitous, and Trustworthy Computing*, volume 2, pages 78–83, 2006.

Intel Mash Maker: Join the Web

Rob Ennals
Intel Research
2150 Shattuck Avenue
Penthouse Suite
Berkeley, CA 94704, USA
robert.ennals@intel.com

Eric Brewer
Intel Research
2150 Shattuck Avenue
Penthouse Suite
Berkeley, CA 94704, USA
eric.a.brewer@intel.com

Minos Garofalakis*
Yahoo Research
2821 Mission College Blvd
Santa Clara, CA 94, USA
minos@yahoo-inc.com

Michael Shadle
Software Solutions Group
Intel Corporation
5200 NE Elam Young Parkway
Hillsboro, OR 97124, USA
michael.shadle@intel.com

Prashant Gandhi
Intel Research
2200 Mission College Blvd
Santa Clara, CA 95054, USA
prashant.gandhi@intel.com

ABSTRACT

Intel[®]Mash Maker is an interactive tool that tracks what the user is doing and tries to infer what information and visualizations they might find useful for their current task. Mash Maker uses structured data from existing web sites to create new “mashed up” interfaces combining information from many sources.

The Intel[®]Mash Maker client is currently implemented as an extension to the FireFox web browser. Mash Maker adds a toolbar to the browser that shows buttons representing enhancements that Mash Maker believes the user might want to apply to the current page. An enhancement might combine the data on the page with data from another source, or visualize data in a new way. Mash Maker is intended to be an integral part of the way the user browses information, rather than being a special tool that a user uses when they want to create mashups.

In order to create mashups from normal websites, Mash Maker must first extract structured data from them. If the web site does not provide RDF data, then Mash Maker extracts structured data from the raw HTML using a community-maintained database of *extractors*, where each extractor describes how to extract structured data from a particular kind of web site.

Categories and Subject Descriptors

H.4.3 [Information Systems]: Information Browsers

General Terms

Management, Design, Human Factors

Keywords

Mashup, Data integration, Personalization, Visualization

1. INTRODUCTION

Historically, the process of writing new queries and creating new graphic interfaces has been something that has been left to the experts. A small set of experts would create applications, and all users would have to make do with what was available, even if it did not quite fit their needs.

Mashups are an attempt to move control over data closer to the user and closer to the point of use. Although mashups are technically similar to the data integration techniques that preceded them, they are philosophically quite different. While data integration has historically been about allowing the expert owners of data to connect their data together in well-planned, well-structured ways, mashups are about allowing arbitrary parties to create applications by repurposing a number of existing data sources, without the creators of that data having to be involved. The importance of mashups is arguably more political and cultural than technical. Mashups are about the “democratization of innovation” [11].

Intel[®]Mash Maker is a project within Intel Research that is aiming to push the mashup envelope a few steps further. Previous work in mashups has followed a model in which a reasonably skilled user uses a special interface to visually compose information from different data sources, creating a new mashed-up application that can then be used by other users. Although this approach has empowered a new class of semi-skilled users to create a vast number of customized applications, specially tailored for particular tasks, or particular groups of people, we believe that the concepts can be taken much further. With Intel[®]Mash Maker, our intention is that normal users should be able to easily create applications and interfaces that are specially customized not only for them, but for the exact task they are performing at that moment. Our aim is mashups for the moment, on demand.

Intel[®]Mash Maker does this by making mashup creation part of the normal browsing process. Rather than having a reasonably skilled user create a mashup in advance as a mashup site that other users browse to, Mash Maker instead creates personalized mashups for the user inside their web browser. Rather than requiring that a user tell a mashup tool what they want to create, Mash Maker instead watches what information the user looks at, correlates the user’s behavior with that of other users, and guesses a mashed up application that the user would find useful, without the user even having to realize they wanted a mashup.

Mash Maker is currently implemented as an extension to the

*Work done primarily while at Intel Research

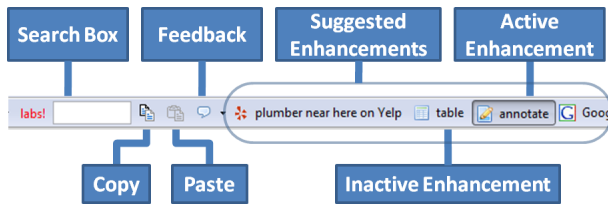


Figure 1: The mashup bar

Firefox web browser, with versions planned for other browsers. As the user browses the web, the Mash Maker toolbar displays buttons representing enhancements that Mash Maker thinks the user might want to apply to their current page (Figure 1). The user need simply turn on some combination of these enhancements to create a new mashup; e.g., to plot all items on a map, the user can click on the Google Maps button.

In addition to suggesting known enhancements defined by other users, Mash Maker will also suggest new enhancements that it has created by filling in *gaps* in known enhancements (Section 4.3). Similarly, a user can create a composite mashup by turning on several generic enhancements (e.g. good restaurants + crime level + map).

If a user knows what enhancement they want and Mash Maker does not suggest it then the user can use a simple copy and paste interface to show Mash Maker a pair of web sites that the user would like Mash Maker to combine (Section 4.2).

If a web site source exports its data in a structured form such as RDF then Mash Maker can use this, otherwise Mash Maker must extract structured data from the raw HTML. Mash Maker consists of two key parts: the client, which is a browser extension that allows a user to create mashups as part of their normal browsing process; and the server, which stores *extractors* that tell Mash Maker how to extract structured data from normal web sites (Figure 2). The server operates like a wiki, allowing any user to edit the extractor for a page.

We believe that Mash Maker offers a radically new approach to querying and visualizing data:

- **Mashups for me, right now, on demand.** Mash Maker allows an unskilled user to create mashups that are tailored not only for them, but tailored for the task they are performing right now.
- **The Mashups come to you.** Mash Maker watches what the user does and tries to suggest mashups that the user will like.
- **Mashing is browsing.** Mash Maker does not require a user to use a special mashup creation interface, or browse to special mashup sites. Instead, Mash Maker augments the familiar web browsing interface that the user already uses to browse data, and enhances this with mashed up information.
- **Rely on the community to structure your unstructured data.** Mash Maker can mash up data from web sites that have no structured data API. It does this by maintaining a centralized community-maintained database of *extractors*. Any user can edit this database to create an *extractor*, describing how to extract structured data from a particular kind of web page.

This paper should be approached as an overview paper. Many of the topics discussed in this paper contain considerable subtleties that we do not have space to explore fully.

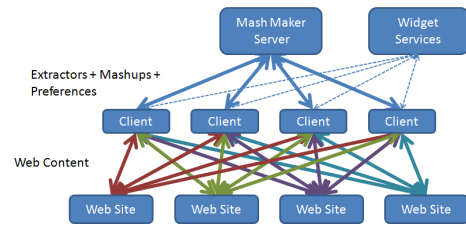


Figure 2: The Mash Maker client and server

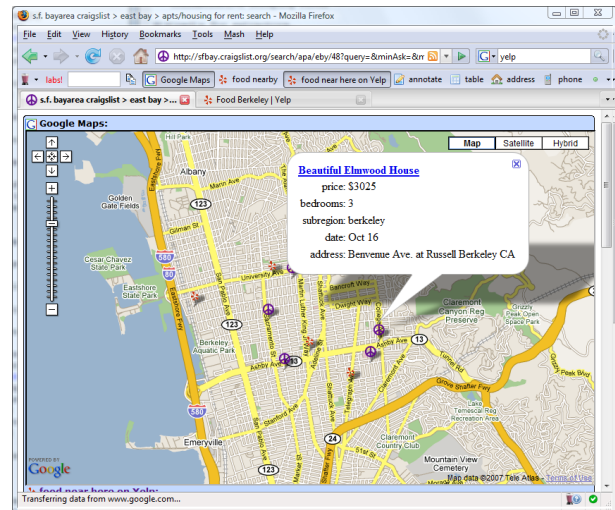


Figure 3: Visualizing data on a map

2. A QUICK TOUR OF MASH MAKER

We will begin with a couple of example usage scenarios for Intel[®] Mash Maker, illustrating some of the concepts that we discuss in this paper:

2.1 Showing things on maps

Over half of the mashups listed on ProgrammableWeb.com involve plotting things on a map. Since this is a common mashup scenario, we will use that as our first example, by showing how a user, Alice, creates the classic “Craigslist houses on a map” mashup using Mash Maker.

Alice browses to the Craigslist apartment listing, as she would normally, and browses the apartments that are available. Mash Maker notices that the page contains items with addresses, and so displays a “Google Maps” button that Alice can use to visualize this data on a map. Alice is interested in an apartment that has good restaurants nearby, so she opens another window and searches for restaurants on Yelp. Mash Maker notices that Alice is interested in restaurants and so updates the mashup bar for the Craigslist page to suggest adding a list of restaurants in the area.

Alice clicks on this new suggestion and Mash Maker responds by inserting information about Yelp restaurants directly into the Craigslist page. Mash Maker found the Yelp restaurants by passing the current location and the topic “food” as arguments to a form on the Yelp website. In effect, Mash Maker has performed a join on the data in Craigslist and Yelp, while obtaining the data through the standard web interface.

Alice also clicks on the “Google Maps” button to see all this data visualized on one map (Figure 3). Finally, Alice turns on the “annotate” enhancement. This is a built-in enhancement that can be suggested for any page that contains items with their own URLs.

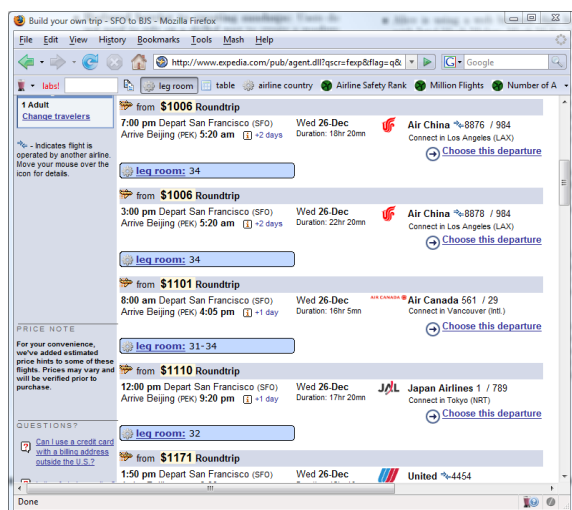


Figure 4: Flights annotated with legroom

The annotate enhancement augments each item on the page (in this case apartments) with an interactive widget that allows Alice to attach a persistent personal note. Items are identified by the (normalized) URLs, and so the same note will appear for each item if it appears on a different page, provided the item's URL is the same. As demonstrated by the “annotate” enhancement, not all enhancements bring in data from other web pages. Some enhancements add new UI features or services that the user might want.

2.2 Information about Flights

As a second example, we will demonstrate the general principle of “copy and paste mashups” by having a user, Bob, add leg room information to Expedia flights.

Bob browses to Expedia and searches for flights. Bob is concerned that he should book a flight that will give him a lot of leg room, but Mash Maker is not currently suggesting legroom. Bob is disappointed that Mash Maker has not given him a button to add leg room to the page, so he browses to a page giving the typical leg room for different airlines. Mash Maker has still not guessed what Bob wants, so he clicks “copy” on the leg room and “paste” on the flight listing, to tell Mash Maker that he would like to augment the page about flights with information from the page about legroom. Mash Maker brings up a dialog box, asking Bob how the data should be connected. Mash Maker has correctly guessed that the data should be equi-joined based on the “Airline” column, but is not sure with which column to annotate the flights. Bob selects “economy legroom”, gives the new enhancement a brief name and description, and clicks “publish” to tell Mash Maker that this enhancement should be suggested to other users. Mash Maker turns the new enhancement on automatically, giving the display in Figure 4.

3. THE BASIC APPROACH

The examples in Section 2 demonstrated several of the features of the Intel[®] Mash Maker user interface. The user can turn on and off any combination of enhancements by clicking on buttons on the mashup toolbar (Figure 1). The user can also search for a specific enhancement by typing keywords into the search box. If Mash Maker does not suggest appropriate enhancements then the user can show Mash Maker what they want by viewing web pages about the topic they are interested in, or by using copy and paste to explicitly tell Mash Maker what they want to combine.

Mash Maker allows the user to pick and choose the information and visualizations that they would like to combine together. Going back to the example from Section 2.1, there are many different apartment listing sites, many different local search services, and many different mapping services. With Mash Maker, a user can pick any combination of these services by simply browsing to their preferred apartment listing service, and then turning on the enhancements for their favorite local search service and mapping service. It is not necessary for any pair of these services to have been combined previously, so long as other users have previously shown Mash Maker the generic ways that these sites can be combined with other sites (Section 4.2).

When enabled, an enhancement adds additional content to the current web page (Figure 4). This information is visualized using a widget, which may be static or interactive. Enhancements are non-side-effecting. Clicking an enhancement button can only add information to the page — it cannot perform externally visible actions. If an enhancement wants to perform actions, then it must do this by inserting a widget that the user can use to perform such actions. For example, rather than having an “add to my calendar” enhancement (as in Operator [18]), one would instead have a “calendar operations” enhancement that adds an “add to my calendar” button to every event on the page.

Internally, Mash Maker describes its enhancements using an underlying functional programming language [5].

Previous mashup tools [8, 19, 20, 21] have taken a server-based approach in which a mashup server retrieves data from other sites and uses this to create a new web site that hosts the mashup. Mash Maker instead runs almost entirely on the client, running as an extension to the user's web browser. This client-based approach has advantages for data access, privacy, performance, and user experience. Since Mash Maker is a browser extension, Mash Maker can see everything the browser can see, including local files, information on the intranet, information requiring a login, and active content generated by Javascript. There are also privacy and performance advantages, since Mash Maker does not need to ship data to and from a central server in order to create the user interface. Finally, running as an integrated part of the browser allows Mash Maker to present a more pervasive user experience, in which mashups are an integral part of the way the user looks at information.

A key driving principle of Mash Maker is that mashups should be personal. Rather than using fixed mashups created by other users, an unskilled user should instead be provided with mashups that have been created specially for them, and for the task that they are currently performing. Mash Maker looks at the information that the user browses, and the mashups that the user has turned on in the past, and uses this to guess what mashups this particular user might want right now, which it then suggests on the mashup bar. Since Mash Maker runs on the client, it can mine sensitive private information without leaking it to third parties. In future versions, we plan to make use of physical sensor information such as location, time, device type, and inferred activity, to improve the suggestions that Mash Maker makes.

3.1 The Server

In order to create such mashups, Mash Maker needs to understand the meaning of web pages. In an ideal world, all web sites would expose the structured databases that underly their sites, making it easy to integrate the data, however this ideal world has not yet arrived. It is thus necessary, at least for now, that we use more ad-hoc techniques to understand the meaning of web pages.

Mash Maker understands the meaning of web pages using a col-

laboratively maintained database of *extractors*. An extractor describes how to extract structured information from the raw HTML of a particular kind of page. For example, the extractor for Craigslist apartments says how to find an apartment on the page, and how to extract each of the properties of an apartment.

The Mash Maker server is influenced by wikis such as Wikipedia. Like a wiki, the Mash Maker server allows any user to edit the extractors for any web site. To avoid vandalism, Mash Maker allows high profile or sensitive pages to be locked down so that they can only be edited by trusted users. Mash Maker also provides a complete version history, allowing users to roll back previous edits if vandalism has occurred. Section 3.4 briefly discusses some of the security issues relating to bad data.

In addition to being able to tell us how to extract meaning from a page, the Mash Maker server also stores information about how the page is parameterised. For example, if we have a page about “England”, then the server can tell us that the page is parameterized by a country, and that the argument is encoded as a form parameter of the URL. A related URI-comprehension mechanism is also used for normalizing URIs that are textually different but refer to the same resource. This information is provided by a collection of *arg handlers*, which are managed similarly to extractors.

The Mash Maker browser extension includes an extractor editor, allowing any user to edit the semantic extractor for the page they are currently browsing by opening the extractor editor side bar. Indeed Mash Maker will prompt the user to do this if it does not understand the meaning of the current page.

3.2 Suggestions

Mash Maker chooses which enhancements to suggest using an ad-hoc algorithm that assigns weights to enhancements based on a number of factors. The main factor affecting the weight of an enhancement is how recently and how often the current user and other users have applied that enhancement to pages similar to the current page. For each extractor/enhancement pair, the Mash Maker server maintains a record of how often and how recently the user and the community as a whole have applied that enhancement to pages described by that extractor.

The suggestion algorithm also uses a number of ad-hoc heuristic rules to improve suggestion weights, including favoring information from sites that the user has viewed recently, and taking account of explicit votes for and against particular enhancements by users. The current suggestion algorithm is quite crude and we believe there is much potential for improvement. We plan to improve it in future work.

3.3 Copyright

We have no interest in using content in ways that the creators disapprove of. However, it is impractical to ask permission from every site in advance, since we don’t know what content our users might wish to combine. Our approach is to assume that a small amount of data extraction from a website is probably harmless. If we see that a data source is being used a lot, then we will contact the owners of that content to ask them if and how they would like their content to be used, and store this information on the central server. For content for which we do not yet have an agreement, we throttle the rate at which Mash Maker extracts data until we know what the content owner wants. Our hope is that, just as most web sites like being listed by Google, most content owners will appreciate the additional exposure Mash Maker provides for their content.

3.4 Privacy and Security

Giving unskilled users the power to combine data sources in previously unexpected ways opens up a number of issues for privacy and security. If one applies an enhancement to a page that contains private information, then that enhancement could cause that information to become visible to a third party. For example the “Google Maps” enhancement sends all addresses on the page to Google, which might not be acceptable if they were the addresses of confidential locations. There are several ways we try to address this problem. First, all enhancements are manually checked by trusted “moderator” users, before they can be suggested to other users, to make sure that they are not obviously malicious. Second, the Mash Maker server has facilities for marking data on a page that should be considered private, and should not be passed outside the client. Third, some standard classes of confidential data, such as passwords, and credit card numbers, can be easily detected and blocked from being passed outside the client. This is an area of active research, and we do not yet have a perfect solution.

A similar issue is the ability of Mash Maker to track what users do. One of the goals of Mash Maker is that the client should tell the server as little as possible about what the user is doing. In particular, when the client requests extractors from the server, it does so for an entire domain rather than a particular page, and does not send identifying information (while we could log IP addresses, we intentionally avoid doing so). Moreover, since the Mash Maker provides extractors using an anonymous, cacheable, REST API, a group of users could potentially hide their behavior from Mash Maker by accessing it through a shared proxy. Since a user may have multiple devices and browsers that they wish to be synchronized, Mash Maker will, by default, store on the server a record of what mashups the user has applied to particular general kinds of page. The user can turn this off if they prefer.

3.5 Query Optimization

Extracting data from web pages is often a very inefficient way to access a data source. Our hope is that the popularity of mashup tools will encourage an increasing number of content providers to provide high level access to their data, through interfaces such as SPARQL [23]. If we know that the data provided by a site is not private, and that it is okay for us to cache it, then the Mash Maker server will cache it in a high-level database on the Mash Maker server. The Mash Maker client can use this database to obtain data using efficient high-level queries. For example, to retrieve the head of state of 100 countries using Wikipedia pages would require downloading 100 pages if accessed directly, or doing a small database query to the cache on the Mash Maker server. We believe that further research is needed on the optimization of cross-provider queries of web content.

4. MAKING MASHUPS

The enhancements suggested by Mash Maker take several forms, from simple linked data, through to mashups inferred from user behavior, and new visualizations. Our experience so far is that, while most users limit themselves to turning on combinations of previously defined enhancements, one only needs a fairly small number of more skilled users to create enhancements for all users to benefit, since all the enhancements these users create can be used by other users.

4.1 Linked Data

The simplest kind of enhancement is one that just follows a link on the page and inserts information that it finds there. If an item on a page contains a URL for another page, then Mash Maker will automatically provide enhancement buttons for annotating the current page with information described on the linked page.

4.2 Copy and Paste

If the user knows what mashup they want, and Mash Maker does not automatically suggest it, then they can teach Mash Maker new connections between web sites using a simple “copy” and “paste” interface. The user clicks “copy” on the *source* page that they would like to use information from, and then clicks “paste” on the *host* page that they would like to add this information to. For example, in Section 2.2 the user copied information about airlines and pasted it into a page listing flights with those airlines. Mash Maker will try to guess how the data should be combined, based on the property and form argument types for the two pages, and the behavior of past users. The most common ways to combine pages are to do a simple join of the data, and/or pass data from the host page as a form argument to the source page. If Mash Maker guesses wrong then the user can manually specify how to combine the data.

Of course, if the user has to edit the enhancement manually then we are essentially back to the same difficulty level as specifying a query in a database. Mash Maker thus tries where possible to avoid the user having to do this, either by guessing how to combine data, or by suggesting enhancements created by previous users.

4.3 Filling in the Gaps

When creating an enhancement using copy and paste, one can leave *gaps* in the definition that can be filled in later. These gaps correspond to function parameters in the underlying functional language [5]. Mash Maker can fill in gaps with information from pages the user looked at recently. For example, in Section 2.1 we used a local search enhancement that allowed Mash Maker to guess what the user was interested in. In that case, Mash Maker guessed that the user wished to search for “food” because Mash Maker had seen this was the “search term” property of a previous page. It is also possible for a user to fill in the gaps explicitly, by clicking on the button for the enhancement template, and entering the arguments directly.

4.4 Applicability of an Enhancement

Once a user has defined a new enhancement, Mash Maker can suggest the enhancement for any pages that are *similar* to the original host page. Early versions of Mash Maker could potentially suggest an enhancement for any page that had the properties that were required by the enhancement, or that had applicable enhancement that could produce such properties. However we found that this caused Mash Maker to suggest a lot of inappropriate enhancements. More recent versions allow a user to restrict the classes of items for which a particular enhancement should be suggested. For example, one might constrain the “legroom” enhancement to only be applicable to “flights”.

4.5 Visualization Widgets

Mash Maker visualizes data added to a page using *visualization widgets*. Widgets range from simple static widgets such as static text and images, through to right interactive widgets such as maps. A user can write a new widget by creating a normal web page that exposes a Javascript function called `mashmaker_widget`. Mash Maker renders a widget by inserting this page as an `iframe` [12] and passing the javascript function a value representing the RDF data

that the widget should visualize. Widgets can communicate with other web services. For example the *annotate* widget communicates with a web service that stores personal notes written by users, and the *google maps* widget communicates with Google Maps.

4.6 Schema Matching

Since content providers are not always consistent in the names they use to refer to the same thing, Mash Maker allows users to interactively teach Mash Maker which strings and URLs should be considered equivalent. If a mashup is trying to join together two data sets and can’t find a connection then it will insert a button saying “click to connect this”. If the user clicks this button then they are presented with an interface that allows them to say what the key should have been matched to. As with all other metadata, such equivalences are edited collaboratively, and are shared with all users.

More generally, combining information from arbitrary web pages, decoded using extractors written by different authors, is a potentially very difficult schema-matching problem. We have so far only touched on the issues that need to be addressed. While our current simple implementation is already able to do a good job in many cases, we intend to explore this area a lot more in the future, including looking at what ideas we can borrow from previous data integration work.

5. RELATED WORK

Since mashups are a hot topic right now, there has been a lot of previous work done in this field.

5.1 Data Integration

The database community has done a huge amount of research on data integration – reliably connecting together data from different sources that might have very different schemas, and different ways of naming the same things. Recent interesting examples include SEMEX [3], DataSpaces [7], and Cohera [24]. We believe that many of the techniques developed by this work are applicable to Mash Maker. Indeed we see Mash Maker, and mashups in general, as being primarily about providing approachable environments that allow arbitrary unskilled users to easily apply existing data integration techniques to repurpose data from arbitrary existing data sources.

5.2 Mashup Creation as a Separate Activity

There are many mashup tools that allow one to create mashups by graphically combining data sources and operators together as graphical dataflow graphs or pipelines. Examples include Yahoo Pipes [19], Marmite [25], Microsoft Popfly [20], IBM QED Wiki [21], and Anthracite [1]. These are all very powerful tools, and there are interesting differences between them, however they all follow a broadly similar model in which a reasonably skilled user creates a mashup by visually connecting components together, with the intention of creating a new site that users can use. Mash Maker extends this work by integrating both the creation and the use of mashups with the normal browsing experience, and predicting what the user wants based on their past behavior. In addition, with the exception of Marmite, these are all server-based tools, which means they have to deal with the issues we discussed in Section 3.

There are also a number of mashup creation tools that work at a lower level. Tools like Google Mashup Editor [8], Plagger.org, Ning.com, Javascript Dataflow Architecture [14], and Web Mashup Scripting Language [22] give the user a lot of power over the

mashups that they create, at the expense of requiring the user to write some form of program.

5.3 Mashups as Browsing

Operator [18] and Miro [6] are browser extensions that suggest actions to be performed on items they have found on the current web page. Operator looks for data tagged with microformats [15] and Miro uses a sophisticated data detector. Miro allows users to teach it new operations using a “program by example” interface. Mash Maker goes beyond what is possible with these tools by adding content to pages, and allowing one to create complex composite mashups that go beyond applying a single operation to the elements on a page.

GreaseMonkey¹ allows users to write scripts that can arbitrarily change the behavior of websites. If a user visits a web page for which they have registered a script, the greasemonkey script will run and can do pretty much anything to the page, including adding information to the page from other sites, or bringing in extra information. Many greasemonkey scripts provide behavior that is equivalent to, or superior to that which is achievable with Mash Maker mashups. GreaseMonkey provides mashup-writers with enormous power, at the cost of requiring them to write their mashups as Javascript programs.

There are also a number of mashup tools that excel in creating a particular kind of mashup. Google MyMaps and Microsoft MapCruncher² make it easy for end users to create mashups involving maps and Swivel.com makes it very easy for end users to create graph mashups from multiple data tables.

5.4 Semantic Web Browsers

Like Mash Maker, semantic web browsers such as Tabulator [2] and PiggyBank [13] are implemented as FireFox extensions, and allow one to browse data that can be found by following links on a page. PiggyBank allows one to create extractors using Solvent, which is similar to our extractor editor. Mash Maker improves on these tools by allowing one to add new information directly to the current page, rather than having to use a new interface. Mash Maker also extends this previous work by allowing users to compute new data from the data they have available, rather than being restricted to finding information by following links.

There are many other semantic web browsers. For example Haystack [9] is implemented as a stand-alone application and OntoWiki [10], DISCO [4], mSpace [16], and OpenLink [17] are implemented as web applications. All semantic web browsers we are aware of treat semantic web browsing as a separate task, rather than augmenting the normal browsing interface with semantic data.

5.5 Widgets

A huge number of sites exist that provide widgets that can be embedded into pages to show visualization of data. Google Widgets, ClearSpring.com, Widsets.com, WidgetBox.com, and Apple’s Dashboard allow users to write small graphical web widgets and then lay them out together on a screen. DataMashups.com additionally allows users to connect these widgets together.

The difference between Mash Maker widgets and these other widget systems is more in their intended purpose, rather than the underlying architecture. Mash Maker intends that widgets be used to visualize potentially arbitrary RDF data, rather than being loosely parameterised representations of a particular web site.

¹<http://diveintogreasemonkey.org>

²<http://research.microsoft.com/mapcruncher>

5.6 Content Suggesters

Many tools exist that try to understand the user’s interests and suggest things that they might want. StumbleUpon.com is a browser toolbar that suggests web sites that one might be interested in. Last.fm and Pandora.com are internet radio stations that try to play songs that the user would be interested in. Amazon.com has a product suggestion system that suggests products that you might be interested in, based on past behavior.

6. CONCLUSIONS

Mash Maker as it is now is just a small step toward our eventual vision. Our intention is to move toward a personal proactive internet in which a user’s computing devices anticipate what the users wants and make use of semantic information on the internet to present them with the information they want, presented the way they want it, while requiring a minimum of interaction.

Mash Maker is currently available as a limited “technology preview release”. Although use is currently invite only, members of the public can sign up online to be put on a waiting list to be sent an invite when we want more testers. Visit the URL below to try out Intel[®] Mash Maker:

<http://mashmaker.intel.com>

Acknowledgments

This work has benefited from the input of many people. Particular thanks should go to Kulki Dattatraya, David Gay, Badari Kommandur, Eric Paulos, Rusty Sears, Ian Smith, and K Sridharan.

7. REFERENCES

- [1] Anthracite. <http://www.metafly.com/products/anthracite>.
- [2] T. Berners-Lee, Y. Chen, L. Chilton, D. Connolly, R. Dhanaraj, J. Hollenbach, A. Lerer, and D. Sheets. Tabulator: Exploring and analysing linked data on the semantic web. In *Proceedings of the 3rd International Semantic Web User Interaction Workshop*, 2006.
- [3] Y. Cai, X. L. Dong, A. Halevy, J. M. Liu, and J. Madhavan. Personal information management with SEMEX. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*, pages 921–923, New York, NY, USA, 2005. ACM Press.
- [4] Disco - hyperdata browser. <http://sites.wiwiiss.fu-berlin.de/suhl/bizer/ng4j/disco/>.
- [5] R. Ennals and D. Gay. User-friendly functional programming for web mashups. In *ICFP '07: Proceedings of the 2007 ACM SIGPLAN international conference on Functional programming*, pages 223–234, New York, NY, USA, 2007. ACM Press.
- [6] A. Faaborg and H. Lieberman. A goal-oriented web browser. In *CHI '06: Proceedings of the SIGCHI conference on Human Factors in computing systems*, pages 751–760, New York, NY, USA, 2006. ACM Press.
- [7] M. Franklin, A. Halevy, and D. Maier. From databases to dataspace: A new abstraction for information management. In *SIGMOD Record*, 2005.
- [8] Google mashup editor. <http://editor.googlemashups.com>.
- [9] Haystack project. <http://groups.csail.mit.edu/haystack/>.
- [10] M. Hepp, D. Bachlechner, and K. Siorpaes. Ontowiki: community-driven ontology engineering and ontology usage

- based on wikis. In *WikiSym '06: Proceedings of the 2006 international symposium on Wikis*, 2006.
- [11] E. V. Hippel. *Democratizing Innovation*. MIT Press, 2006.
 - [12] HTML 4.01 specification.
<http://www.w3.org/TR/REC-html40/>.
 - [13] D. Huynh, S. Mazzocchi, and D. Karger. Piggy bank: Experience the semantic web inside your browser. In *Proceedings of the 4th International Semantic Web Conference*, 2005.
 - [14] S. C. S. Lim and P. Lucas. Jda: a step towards large-scale reuse on the web. In *OOPSLA '06: Companion to the 21st ACM SIGPLAN conference on Object-oriented programming systems, languages, and applications*, pages 586–601, New York, NY, USA, 2006. ACM Press.
 - [15] Microformats. <http://microformats.org>.
 - [16] mspace. <http://mspace.fm>.
 - [17] OpenLink RDF Browser. <http://demo.openlinksw.com/DAV/JS/rdfbrowser/index.html>.
 - [18] Introducing operator. <http://labs.mozilla.com/2006/12/introducing-operator>.
 - [19] Yahoo Pipes. <http://pipes.yahoo.com>.
 - [20] Microsoft popfly. <http://popfly.com>.
 - [21] Qedwiki. <http://services.alphaworks.ibm.com/qedwiki/>.
 - [22] M. Sabbouh, J. Higginson, D. Gagne, and S. Semy. Web mashup scripting language (poster). In *16th International World Wide Web Conference*, 2007.
 - [23] SPARQL Query Language for RDF.
<http://www.w3.org/TR/rdf-sparql-query/>.
 - [24] M. Stonebraker and J. M. Hellerstein. Content integration for e-business. In *Proceedings of the 2001 ACM SIGMOD international conference on Management of data*, 2001.
 - [25] J. Wong and J. Hong. Marmite: end-user programming for the web. In *CHI '06: CHI '06 extended abstracts on Human factors in computing systems*, pages 1541–1546, New York, NY, USA, 2006. ACM Press.

Ricardo Baeza-Yates Speaks Out **on CS Research in Latin America, His Multi-continent Commute for** **Yahoo!, How to Get Real Data in Academia, and Web Mining**

by Marianne Winslett



Ricardo Baeza-Yates
<http://www.dcc.uchile.cl/~rbaeza/>

Welcome to this installment of ACM SIGMOD Record's series of interviews with distinguished members of the database community. I'm Marianne Winslett, and today we are at the ICDE 2007 conference in Istanbul. I have here with me Ricardo Baeza-Yates, who is the vice president of Yahoo! Research in Europe and Latin America. Before that, he was a professor of computer science at the University of Chile for many years. His research interests include information retrieval, databases, algorithms, and user interfaces, and he is a co-author of one of the most widely used books on information retrieval. Ric's PhD is from the University of Waterloo. So, Ric, welcome!

You have moved from an academic environment in Chile to an industrial lab in Barcelona---isn't that a double culture shock?

Not really. Barcelona is part of Spain, and Spain was the original power that conquered Chile, so I think culturally it is very similar. People do speak a different language in Barcelona, Catalan, which I can understand but I haven't mastered. The academic to industry transition was not so difficult, because we were already doing almost the same research as we are doing now at Yahoo!.

I have heard of telecommuting, but not between continents. How do you handle the extreme distances between your Yahoo! outposts?

I travel about three or four times every year to Chile, and I am there for around two months each year. Almost all the people we have there are people who have worked with me before, so I really can work with them remotely. Currently, it is a nice challenge to have a real distributed lab, because it does not matter from where you are connected.

So you view your part of Yahoo! as a single lab that is distributed, rather than separate ones on each continent?

Right, I view it as a single lab. In some projects we have only one location, but in other projects we have both locations working on it.

Will you grow in both locations?

Yes, I think so. That is my goal.

What is it like to do IR and DB research in South America?

When I finished my PhD, I wanted to go back and try to make a difference in my country. There was a small group of PhDs trying to build a good computer science department, and I helped to contribute to that. If you look at South America, after Brazil, Chile is the main research power there. In Chile we cannot compare with Brazil in the size and quantity of computer science research, but I think we can try to compare in quality. And now we have a PhD program; we have students coming from most neighboring countries, and even Mexico, so I think we are doing well.

It is not easy because we have to find many sources of funding, because we don't have big funding agencies. But five years ago, the Chilean government started a big program where you can receive a reasonable size of money per year. Then I started the Center for Web Research, which I think is well known now in the world. But of course, everything is more difficult than in the developed countries.

Does most of your money for your university lab come from industry, then?

No, most of the money comes from public funding from the government. Industry doesn't do research in Chile, and even in some developed countries, like Spain, industry doesn't do too much research.

Do you think that is likely to change?

I don't know. I would like that to change. There is a dilemma between buying technology and producing technology. In some fields it is really very expensive to produce technology. In our field, where it is basically ideas and software, it is not so expensive. But sometimes politicians don't know that potential.

What do you see as the future of database research in Latin America?

Right now I think it is confined to what I mentioned earlier. Brazil has a long tradition of database research, and many groups in that area. They have a very strong program for PhDs inside the country, so now they have more than a thousand researchers in computer science. They have a very strong database conference (SBBDB), which is international.

In Chile we are trying to build a very good group. We have a small group at the Center for Web Research, people that studied in, say, Toronto, in the US, and so on. In the rest of South America, the problem is not that there are no good people, but that there are no opportunities; so they usually leave and work elsewhere. If you consider all of Latin America, Mexico is the second power in computer science after Brazil. They are doing well, and they also have a very nice tradition in database research. And the next ICDE will be in Cancun!

Yahoo!'s web site says your TodoCL ('All of Chile') search engine had 2 million users a month, "a rich source of data for Baeza-Yates on how people use a search engine. 'But at Yahoo!, I'll have access to far more data,' Baeza-Yates says. 'You go from gigabytes to terabytes of data.' " While I am happy for you, I see this as a real problem for academia. How can academia retain its leading database researchers? Aren't they all going to be sucked into industry by the siren song of unlimited access to real data?

That's a possibility, but I don't think that will happen. Not all people want to work in industry. For me, it was a very tough decision to move to industry, but I don't regret it now. It is much better than I expected.

Let me tell you about the example of TodoCL. I built a search engine that had a lot of people using it. I did that with very low funding. It was not easy, but having a search engine basically means having a few PCs in a data center, with good software that grinds along. From that search engine we got data that no one else had, apart from the big search engines. That data allowed us to do things that other people couldn't do.

So I don't really believe that academics cannot get access to real data. I believe that academics can be more entrepreneurial in their work, by setting up some services that become used a lot. The main examples in our field are Citeseer and DBLP, which are services that a lot of people use. And DBLP has basically one person behind it, Michael Ley! So I think that if you have the leadership, and you have the will to do it, you can create interesting services and gather a lot of data on how people use them.

My summary of what you just said is that academicians should build their own interesting little programs and collect their own data. Do people still use TodoCL today, or is there a Google version that is specific to Chile that has become popular?

Google was popular even when I started. People use both Google and TodoCL because they give different results. Vertical search engines have the advantage that you can crawl the country much deeper, have more pages, and also have local rankings. You get different search results using Google, Yahoo!, and TodoCL, and all of them are good.

And now the reverse question: won't academic research become increasingly irrelevant for the web, since researchers outside of industrial labs have little to no access to large-scale information about how people behave on the web?

I think I answered that partially. At Yahoo!, we want to do open research, to share our results. We have been sharing data with some people, such as through a special Yahoo! Research Alliance program where you can get access to data. So, we want to have the help of the academic community, and we want to contribute back to the academic community. We don't believe in "black boxes" where you don't know the level of your impact. You have to go to conferences to know where you are. We also don't believe that we have all the answers, that we know all the trends.

I am not familiar with Yahoo!'s data sharing program, but I know the one at Microsoft is very small. Maybe 10-20 researchers get access to Microsoft's Live Labs data. What's the scale of the data sharing program at Yahoo!?

I think right now it is a similar size, maybe a bit larger. But we are trying to increase that.

You have said that the goal of your new Yahoo! research lab is “to shape the future of the internet”. Aren’t a lot of other people trying to do the same thing right now?

I think that is a goal of all of Yahoo! Research, not only my lab. Yes, Google and Microsoft and other companies are trying to shape the future of the internet, but I think the potential for the internet is so large that it has space for everyone. Even in the market space, there is space for everyone.

What do you think of the current confluence of IR and database research? Did you see it coming?

I think I saw it coming. I did work on structured IR before XML became popular. I have one of the seminal papers on that topic, and that was written in 1995, which was early with respect to the web. We are at a point of confluence on the integration of databases and IR. We will be able to understand each other better if we do more of what I did yesterday at ICDE, i.e., having invited talks by IR people in database conferences, and the other way around. We have different views of the same problem, and both are valid. Today, there is much more unstructured information than structured information. On the other hand, with structured information you know much more, so you can do more complex things. There are both pros and cons on every side.

How should we rank results retrieved from different kinds of sources, such as blogs, shopping centers, structured records, images, and ordinary web pages?

I think that is one of the main open problems in the web: how to do good ranking, especially when you have different sources of data and evidence of quality---such as how people use it, how people created it, how people comment on it, and so on.

For example, you mentioned blogs. A blog is usually written by one person, and commented on by other people. One problem behind the Web 2.0 is how to rank *people* rather than data; if we can know that the author is good and you can trust his or her opinion, we can do better ranking for every kind of data. In ranking, we are working to be able to combine information on people, data, and usage. We call that *community systems*, and at Yahoo! that work is being led by a well known person in the database community, Raghu Ramakrishnan, who was at the University of Wisconsin before coming to Yahoo!.

How will you come up with a ranking system that can’t be fooled?

That is very tough. I think today we have to live with a lot of spamming on the web: content spam, link spam, usage spam. We have some answers, but this is a permanent fight. We improve our techniques to detect spammers, and they improve their techniques to deceive us. I hope eventually we’ll end with something closer to a semantic web, where we know whether we can or cannot trust a source of information.

At Yahoo!, you have been seeing scale-up issues for ad hoc queries on thousands of processors, even though the workload is almost embarrassingly parallel. What is the cause of the scale-up problems?

The query processing part is really a parallel problem, so that’s not an issue. The indexing does not parallelize well. You either subdivide the document collection into pieces with separate indexes, or you build one single index and subdivide it, which would take a lot of time to build and maintain. So you have to use some kind of divide and conquer approach in the collection.

But still that means that you need to have very large memory caches just to hold the index, and that is not necessarily only a parallel problem.

What kind of solution do you foresee for that?

I see a real distributed search engine where you have collections divided by say, language, culture, geography. You can also use the network topology, and the query distribution of the region, and so on. You can use a lot of levels of caching, to try to come up with solutions that will amortize the network latency on the internet.

From your experience as a member of the Chilean Academy of Sciences, can you tell us how computer science is viewed by scientists from other disciplines?

I can tell you how other disciplines view us in Chile. I think that view is shared by many places.

Sometimes people believe that computer science is more like a technology that you can use, more like a tool. They don't believe it is a science; they think that it is closer to engineering, but not so important. They think that you can buy a computer, use it, and that is it; they don't see the complexity behind the technology that created that computer and its software. The few who understand us may be people from electrical engineering and mathematics, who work in areas that are closer to our problems.

I think that in time, people will start to recognize that computer science is at the interface of engineering and science. We are not a pure science like physics, and we are not pure engineering like, say, mechanical engineering. That's what makes computer science problems so difficult.

I wrote a long manifesto on this topic in Spanish, and included a joke that says, "The name *computer science* has two problems. The first one is that the name includes *computer*, which is a machine, an object, not an abstract concept. The second problem is that the name includes *science*; if we are a science, we don't need to say that we are a science." It is as though we are too young, and we have to use this name to be sure of our personality or our ideas.

Like library science or political science.

Right, but in this case we are closer to *car science*. With *political science*, at least *political* is not an object, it is an abstraction. The equivalent name in our field would be *computing science*, and some departments do use that name. Even within the field there is disagreement; some departments are called *Computing Science and Software Engineering*. I thought that software engineering was part of computer science, but maybe some people don't.

I also like the quote that you had in that article that said, "In theory there is no difference between theory and practice, but in practice, there is a difference between them."

That quote is from Donald Knuth in his 1999 invited talk at the IFIP World Congress. Many times theoreticians don't turn their ideas into programs, so implementing their ideas looks easy to them. But when you turn a theoretical idea into a program, you may find that the theory was correct but the constant factor in the running time analysis is, say, 10,000, and the approach is not competitive. I think the best theories are inspired by practice and the best practice is inspired by theory, and that is another quote by Don Knuth.

One of my former students had a T-shirt that said, “0 and 1: How Hard Can It Be?” In that vein, please tell us about Al-Khorezm, el matemático olvidado.

That means ‘the forgotten mathematician.’ He was a mathematician from the Caliphate of Baghdad in the Arabic empire. He was so famous that he was called Al-Khorezm, which means the person from Khorezm, the city where he was born. This person is so important to us because he invented the first algorithm, in some sense. (There were previous algorithms by Euclid, Eratosthenes, and other Greeks, but they did not realize what they were doing.) Al-Khorezm was really thinking in logical steps, like a computer algorithm.

Al-Khorezm did so many things. He brought the concept of zero and decimal notation from India. Before that, when there was a zero people wrote nothing, so it was very hard to recognize that the “nothing” was there. Al-Khorezm also wrote the first book on algebra; his book was later translated in Toledo and then reached the rest of Europe. He was the first person who tried to write mathematics for everybody. Al-Khorezm also collected all the works of the Greeks. Thanks to him, those works were not lost in the Middle Ages in Europe.

I think we owe Al-Khorezm almost everything, starting with a very important word in computer science, *algorithm*. The word *Al-Khorezm* is the root of *algorithm* and *logarithm*.

I understand that maps are your hobby. What does that mean?

I like historical maps, so I collect old maps. I like geography, so I like to know places, I like to go to new countries, to meet people, to learn. I like to learn every time I am researching, and not only in computer science. If I see something new in food, I like to try it. I am not afraid of learning new things. Maybe computer science research is one way to realize my hobby, as I have the opportunity to meet many different people, know many different cultures, and visit many different countries. I like what I do. I am thankful for that.

Do you have any words of advice for fledgling or midcareer database researchers or practitioners?

I am not a real database researcher, but I think that there are many interesting problems in the web. From the web we can capture a lot of different kinds of relationships: between content, between structure, between links, between usages of the data. I am sure there is a large potential to do data mining there. So data mining should be one of the main research topics today, especially web mining. And that implies new interesting fields like graph mining, where I think a lot of things can be done. For example, how can we scale up graph mining techniques to mine a graph of one billion nodes? I think that is a very interesting problem. We can use mining to find interesting patterns, things that may imply, say, a new service or something else that people want, or maybe detecting spam. There are many possible applications there, so I think the web has a lot of room for new research problems. That is why the focus of the Barcelona/Santiago lab is Web Mining, which in the end is the Wisdom of the Crowds.

Among all your past research, what is your favorite piece of work?

It is difficult to pick a single one, but I think I have to mention my book with Berthier Ribeiro-Neto, *Modern Information Retrieval*, because it is used all over the world. It has been translated to two other languages. I have been surprised that in many places they know it, and even in small places they use it. It is really amazing how you can reach so many people.

On the research side, I think my initial work on algorithms started a couple of different subfields. I wrote a paper about how to use the parallelism at the bit level in the CPU to speed up string matching with possible mismatches (“A New Approach to Text Searching”, *Communications of the ACM*, October 1992, with G. H. Gonnet). One of my first papers was about a cow trying to cross a fence to go to the other side to eat more grass (“Searching in the Plane”, *Information and Computation*, October 1993, with J. Culberson and G. Rawlins). At the time we did the research (1987), we didn’t know that we were working on one of the first online algorithms for computational geometry. That paper inspired a lot of work on how to find things without having the whole picture, without having the map, so to speak. I think that is a very interesting problem. Those two papers are cited by hundreds of other papers.

If you magically had enough extra time to do one additional thing at work that you are not doing now, what would it be?

I think it would be to stay more with the people I love. I sometimes feel that I don’t have enough time to be with them and also to do the work I like to do. It is very hard when you do work that you like, because it takes time from other things.

If you could change one thing about yourself as a computer science researcher, what would it be?

That is a very hard question. I think, if I could change something, I would like to be better at words. I am not a person who talks too much.

But when you write, you write very eloquently.

Yes, but I am talking about *talking*. Not about *writing*. For me it is much easier to write or to listen than to talk. I think I would like to improve my oratory skills. Some people can stand in front of others and talk about almost anything without having to think about it, and I would like to have that ability.

Thank you very much for talking with us today.

Thank you for the invitation.

Data and Web Management Research at Politecnico di Milano

Stefano Ceri, Cristiana Bolchini, Daniele Braga, Marco Brambilla, Alessandro Campi, Sara Comai, Piero Fraternali, Pier Luca Lanzi, Marco Masseroli, Maristella Matera, Mauro Negri, Giuseppe Pelagatti, Giuseppe Pozzi, Elisa Quintarelli, Fabio A. Schreiber, Letizia Tanca

Dipartimento di Elettronica e Informazione (DEI), Politecnico di Milano
Piazza L. da Vinci, 32 – 20133 Milano, Italy
first_name.last_name@polimi.it

1. INTRODUCTION

Research in data management at Politecnico di Milano has a long and solid tradition; forefront books on distributed databases, conceptual database design, logical databases, and active databases contributed to shape the foundations of this discipline in the last two decades. Historically, our work has addressed both all aspects of innovation in the technology of modern data management systems and the consequent support of design methods and tools.

Recently, a large fraction of the group's efforts has been dedicated to the Web, considered as the ubiquitous infrastructure for making access to distributed and heterogeneous data sources. Our work for the Web concentrated on the models, methods, languages, and tools for supporting the design and automatic generation of modern, data-intensive Web applications.

Although classifying our work in data-driven vs Web-driven is rather arbitrary – and to some extent misleading, because in many projects we use both technologies – for ease of organization we will use this classification. Our report focuses on the last five years of activity.

2. DATA-DRIVEN RESEARCH

In recent years, work addressed the enhancement of database technology in several directions, including active, temporal, spatial, and mobile/very small databases; we also focused on query and mining languages for XML repositories and on supporting effective usage of genomic information. In these fields, emphasis was placed on formally defining the new features required by each specific data management extension, and then inferring properties descending from those definitions, which lead to improved system implementations or to better understanding of the system behaviour. While

supporting the technological advances in data management is important, it is perhaps even more important to support developers and users in taking full advantages of the technology; therefore, database research in the group has always been characterized by an emphasis on innovative languages, methods, design support environments and tools, which could bridge technology to its use in real-life applications.

2.1 Active Databases

Research in active databases has been active since fifteen years, with very visible results worldwide. Recent work focused on defining the formal properties of active rules, enabling the development of more powerful rule analysis tools, and on defining strategies for improving the performances of rule sets under a characterization of the load due to passive and active computations. Active rule analysis is important for rule usability: thanks to rule analysis, rule interaction can be tested and properties such as termination and confluence can be proved. In [1], techniques descending from logic programming are proposed as a new paradigm for rule verification. In [2], the new property of event-trace independence is defined; it guarantees that rule executions are indistinguishable even when we consider an arbitrary sequence of their triggering events. In [3], the scheduling of detached rules is optimized under a characterization of the load due to passive and active computations.

2.2 Temporal Databases

Work in temporal databases focused on the introduction of process modelling aspects within temporal information management, and specifically on the impact of conceptual aspects (temporal constraints for process modelling exceptions that may occur during process execution) upon architectural

issues and choices [4]¹. Work [5] analyzes the potential applications of process modelling and temporal databases to medicine.

2.3 Context-Aware And Mobile Databases

Work on context-aware and mobile databases initially focused on data structures and access methods for improving the performance of data management on small, mobile devices. New efficient logical and physical data structures have been defined for DBMSs running on very small, portable, mobile devices [6]; performance, power consumption, and endurance parameters were optimized using an EEPROM Flash device as storage medium. Starting from the needs of the mobile scenario, we developed a comprehensive methodological framework for integrating, tailoring and delivering context-aware data [7]. The method, supported by a prototype tool – *Context-ADDICT*, demonstrated in [8] –, can be seamlessly applied to large information bases, in order to provide users and applications with the appropriate share of data, tailored to the current context [9]. The research on very small devices has recently drifted towards data management in embedded and pervasive systems, generating a research line on query languages for Wireless Sensors Networks [10] and an SQL-like language for pervasive system.

2.4 Spatial Data

Work on spatial data focused on solving the interoperability problems encountered in building a spatial data infrastructure (SDI), consistent with the INSPIRE directive, and in the development of an integrated interoperability architecture capable of dealing with the semantic mapping and geometric harmonization issues raised by the design of a strongly integrated SDI at regional level. The main achievement in this field has been the development of the GeoUML Spatial Conceptual Model [10-13] and its application in the development of a regional geographic database for “Regione Lombardia”. The model has been adopted by the national committee for standards in geographic data (equivalent to FGDC in USA) at CNIPA (equivalent to NIST in USA), <http://www.cnipa.gov.it/>.

2.5 Query Language Design

The contributions to query language design were focused upon XML and how to make XML repositories more usable, both in terms of user

interface and of retrieval success. We designed XQBE (XQuery By Example) [14], a visual query language using examples of XML as a paradigm for querying XML repositories. XQBE is inspired by QBE invented by Moshe Zloof at IBM Watson Research and available on many products (e.g. MS Access). XQBE allows one to formulate simple queries on top of XML repositories, by drawing annotated trees; the language is formally defined and tools map XQBE to XQuery and XPath. The XQBE environment is referenced from the W3C site linking to XML query language implementations, and has been internationally used as a pedagogical tool for learning XQuery. XQBE is also inspiring recent joint research with IBM Almaden Research Center for enhancing the Clio System, an XML mapping research prototype already partially used by commercial products, by adding object-oriented concepts to it [15]. We also addressed the characterization of graph-based queries (for XML and for temporal databases) by means of model-checking based techniques [16].

2.6 Data Mining

Work on data mining focused on new paradigms, algorithms and execution environments for extracting association rules and sequential patterns from XML repositories, thus enabling classical mining operations for a new and important class of repositories. The research on mining XML repositories has given rise to the development of a rich tool environment, named XMINE, supporting several data mining patterns. The main XMINE operator, described in [17], is based on XPath; it can express complex mining tasks, by indifferently (and simultaneously) targeting both the content and the structure of the data.

Another data mining approach consists in the recognition of frequent patterns within XML documents, and on the use of such patterns as summarized representations of the data; these patterns can then be stored and queried, either when fast (and approximate) answers are required, or when the actual dataset is not available, e.g. it is currently unreachable [20].

Additionally, we worked on extracting unexpected patterns (*pseudo-constraints*) from relational databases. This method reveals properties on database states not declared as constraints, but whose violation instances are interesting facts, hence it considers data mining from a new, fully original perspective [19]. Finally, a complete survey of Web Usage Mining is presented in [18], which surveys about 200 papers published in this area between 2001 and 2005.

¹ The paper [4] had the highest number of downloads from the ACM digital library for a period of 28 months.

2.7 Genomic Data Management

Work on genomic data management has produced methodologies and algorithms to effectively use and mine genomic information in heterogeneous and distributed genomic databases. The work also generated a Web-enabled system, named GFINDER: Genome Function INtegrated Discoverer (<http://www.bioinformatics.polimi.it/GFINDER/>) [21-23], allowing scientists to select and evaluate efficiently and dynamically the most relevant functional and phenotypic information supporting knowledge discoveries in different biomedical experiments. It is a system for discovering, using, and mining a large amount of genomic information and knowledge retrieved from many heterogeneous and distributed databases accessible via the Internet for supporting the evaluation and biomedical interpretation of high-throughput biomolecular experiments. It has been actively used by international Research Centers and Universities: at the time of writing GFINDER Web site received nearly 97,000 accesses by more than 5,500 distinct IPs since its opening in 2004. We also developed, in collaboration with the National Institute of Health, Bethesda (USA), a novel heuristic strategy to filter semantic relations extracted from the scientific literature by using natural language processing [24]. The method allows extracting the valuable genomic functional information with enough quality for subsequent applications aimed at uncovering new biomedical knowledge.

3. WEB-DRIVEN RESEARCH

The growth of Web applications as the fundamental infrastructure for business and social activities has generated a strong interest in methods, environments, and tools supporting their design and deployment. The continuous evolution of technologies calls for a foundational, technology-independent approach, rooted in the tradition of information modelling. The main focus of this research is centred upon a conceptual modelling language, called **WebML** (Web Modelling Language)². WebML describes a conceptual model of the Web application in which the various aspects of the specification (respectively

² **WebML** is extensively used in research and teaching by about 50 national and international institutions, including Technion Haifa (Israel), ETH Zurich (Switzerland), Katholieke Universiteit Leuven (Belgium), University of California San Diego (USA), TIFR (India).

WebML was also independently extended by other research groups, as demonstrated by several research papers published at WWW, ICWE, ECBS, SEKE.

the content, the hypertext, and the presentation) are orthogonally combined. The most innovative aspect of WebML is the modelling of hypertexts as collections of elementary units and links, where the units describe both the visualization of elementary elements of a Web page and the operations performed by the application, and links between units capture the user behaviour. The model was initially focused on the display and management of contents, but it has been progressively extended to incorporate other features of modern Web applications, including: process management, Web service invocation and publication, management of adaptive and reactive computations; management of collaborative applications; methods for improving the accessibility, and more in general the quality of Web applications; support for new media, technologies and architectures, including rich Internet applications, VOIP, and Web architectures for embedded systems. The WebML model and design method are patented in US [25] and described by an international book [26].

A WebML specification is a graph, therefore WebML specifications are supported by a visual design tool with extensible components; such tool, called *WebRatio*, has been initially developed at Politecnico as result of EU-funded projects in the fourth and fifth framework, then has been the core of the spin-off company *Web Models*³.

Recent work in WebML addresses Web services invocation and publication, process management, model-driven design and translation, support of semantic Web services, development of rich Internet and of embedded Web applications.

3.1 Service Invocation and Publication

Web service invocation and publication is explored in [27], where WebML is extended to model complex interactions with Web services. The concepts presented in this paper were fully implemented, as described in [28]. Industrial applications are described in the joint work between our group and the spin-off Web Models [29].

3.2 Process Management

Process management was addressed in [30], by extending the WebML modelling language to describe process-enabled Web applications, i.e. applications where the navigation of the user is

³ **Web Models** is a spin-off company participated by the Politecnico having now about 20 employees, subdivided in the two locations of Milano – for commercial development – and Como – the software factory (<http://www.webratio.com/>).

driven by workflow constraints. The process modelling phase drives the hypertext generation phase, by automatically generating (low-level) hypertext skeletons from (high-level) process models, according to different styles of process enactment. In the above context, special emphasis is given to reverse engineering, i.e. the ability to reconstruct the process from the generated hypertext when the latter is subject to modifications and evolution. This work represents one of the first Web engineering proposals for modelling data- and process-centric applications. Due to the loose control on Web clients, exceptions occurring during execution of processes on the Web are a hard problem. Papers [31,32] present a high-level model enabling case-based exception resolution.

3.3 Model-Driven Design

Model-driven design has been the common driver in many research efforts. We have concentrated on collaborative applications in [33], on context awareness in [34] (with an application to the e-learning context in [35]), and on application quality in [36]; quality depends on conceptual properties of the designed applications instead of interface-specific properties. The method is supported by a tool automating the evaluation process. In [37] we address a general method for designing Web applications which uses the new notion of “Web mart”, an extension of the notion of data mart which is suited to Web applications.

3.4 Model Transformations

Model transformations have adapted WebML to the Model-Driven Architecture (MDA) context. WebML has been generalized using UML Meta-Object Facility (MOF), to make it consistent with OMG’s MDA. Code generation techniques are generalized into an abstract model-transformation framework, capable of addressing such tasks as: the generation of metric models for evaluating different size measures of a project (e.g., for automatically producing the Function Point count from a conceptual application model); or the generation of models for driving the automatic testing of applications.

3.5 Support of Semantic Web Services

Support of Semantic Web Services (SWS) *is* fundamental for spreading SWS. A joint team with CEFRIEL⁴ has produced SWEET, a WebML-based

environment for designing applications of SWS that automatically derives a large portion of SWS annotations from their high level models [38], thereby reducing the efforts required by experts. The environment was experimented for developing Web Service Meta Object (WSMO) components (services and mediators), and has participated to the Semantic Web Challenge [39], designed by Prof. Charles Petrie (Stanford University), which took place at Stanford, Budva, and Athens in 2006, and at Innsbruck and Stanford in 2007. This research received an IBM Faculty Award in 2006.

3.6 RIA and Embedded Applications

WebML conceptual model and WebRatio code generation technology were extended along the direction of rich Internet applications (RIA) in order to transfer more application logic from the server to the client. In addition, Web architectures and applications were downscaled to embedded systems, adapting the conceptual modelling primitives, runtime architectures, and code generation techniques to the space and time constraints of embedded architectures. Embedded applications are envisioned in domotics, intelligent buildings, cultural heritage, and industrial automation.

4. CONCLUSIONS

In October 2007, we held a one-day workshop dedicated to new research directions. We gathered about 40 people, including professors, researchers, and students from DEI, CEFRIEL and Web Models. We agreed that emphasis in the future will be dedicated to five major themes: bio & nano technologies, mediation & mapping, social analytics for the Web, new user experiences, and stream reasoning. The last topic, which is not covered in this paper, addresses the massive computation of reasoning (e.g., logical rules) on streaming data, and will be covered within FP7 by a joint research unit DEI-CEFRIEL. We are discussing each theme within a Wiki, and we will be glad to grant read access to anyone interested to contribute.

⁴ CEFRIEL is a Center for ICT excellence, set up in 1988 as a consortium whose components are Academia (represented by Politecnico di Milano, Università degli Studi di Milano and Università degli Studi di Milano – Bicocca), Enterprises (including

some of the most important ICT companies operating in Italy), and Public Administration, represented by the Lombardy Region. <http://www.cefriel.it/>.

5. REFERENCES

- [1] S. Comai and L. Tanca. Termination and confluence by rule prioritization. *IEEE Transactions on Knowledge and Data Engineering*, 15(2):257-270, 2003.
- [2] A. Bonifati, S. Ceri, and S. Paraboschi. Event trace independence of active behaviour. *Information Processing Letters*, 94(2):71-77, 2005.
- [3] S. Ceri, S. Paraboschi, G. Serazzi, and C. Gennaro. Effective scheduling of detached rules in active databases. *IEEE Transactions on Knowledge and Data Engineering*, 15(1):2-13, 2003.
- [4] C. Combi and G. Pozzi. Architectures for a temporal workflow management system. In Proc. of the 2004 ACM Symposium on Applied Computing (SAC), Nicosia, Cyprus, March 14-17, 2004. ACM Press, New York, NY, 659-666.
- [5] K.P. Adlassnig, C. Combi, A.K. Das, E.T. Keravnou, and G. Pozzi. Temporal representation and reasoning in medicine: Research directions and challenges. *Artificial Intelligence in Medicine*, 38(2):101-113, 2006.
- [6] C. Bolchini, F. Salice, F.A. Schreiber, and L. Tanca. Logical and physical design issues for smart card databases. *ACM Transactions on Information Systems (TOIS) ACM*, 21(3):254-285, 2003.
- [7] C. Bolchini, C. Curino, F.A. Schreiber, and L. Tanca. Context integration for mobile data tailoring. In Proc. IEEE Int. Conf. on Mobile Data Management (MDM), Nara (Japan), 2006, 5.1-5.8.
- [8] C. Bolchini, C.A. Curino, E. Quintarelli, F.A. Schreiber, and L. Tanca. Context-ADDICT: a tool for context modeling and data tailoring. In Proc. IEEE Int. Conf. on Mobile Data Management (MDM), (demo paper), May 2007.
- [9] C. Bolchini, C.A. Curino, G. Orsi, E. Quintarelli, R. Rossato, F.A. Schreiber, and L. Tanca. And what can context do for data? *Communications of the ACM*, (in press).
- [10] Schreiber F.A. Automatic generation of sensor queries in a WSN for environmental monitoring. In B. Van de Walle, P. Burghardt, and C. Nieuwenhuis, eds. Proc. 4th Int. ISCRAM Conference, Delft, May 2007, 245-254.
- [11] A. Belussi, M. Negri, and G. Pelagatti. Modelling spatial whole-part relationships using an ISO TC 211 conformant approach. *Information and Software Technology*, 48: 1095-1103, 2006.
- [12] A. Belussi, M. Negri, and G. Pelagatti. An ISO TC 211 conformant approach to model spatial integrity constraints in the conceptual design of geographical databases. In Advances in conceptual modeling theory and practice. LNCS 4231, Springer, 2006, pp.100-109.
- [13] A. Belussi, M.A. Brovelli, M. Negri, G. Pelagatti, and F. Sansò. Dealing with multiple accuracy levels in spatial databases with continuous update. 7th International Symposium on Spatial Accuracy Assessment in Natural Resources and Environmental Sciences, Lisbon, Portugal, 5-7 July 2006.
- [14] D. Braga, A. Campi, and S. Ceri. XQBE (XQuery By Example): a visual interface to the standard XML query language. *ACM Transactions on Database Systems*, 30(2):398-443, 2005.
- [15] A. Raffio, D. Braga, S. Ceri, P. Papotti, and M.A. Hernandez. Clip: a visual language for explicit schema mapping. In Proc. IEEE Int. Conf. on Data Engineering, 2008 (in press).
- [16] E. Quintarelli. Model-checking based data retrieval: an application to semistructured and temporal data. LNCS 2917, Springer Verlag, 2004.
- [17] D. Braga, A. Campi, S. Ceri, P.L. Lanzi, and M. Klemettinen. Discovering interesting information in XML data with association rules. ACM-SAC 2003, Melbourne, USA, March 2003, pp. 1163-1167.
- [18] E. Baralis, P. Garza, E. Quintarelli, and L. Tanca. Answering XML queries by means of data summaries. *ACM Transaction on Information System*, 25(3), 2007.
- [19] S. Ceri, F. Di Giunta, and P.L. Lanzi. Mining constraints violations. *ACM Transactions on Database Systems*, 32(1):6, 1-32, 2007.
- [20] F.M. Facca and P.L. Lanzi. Mining interesting knowledge from Weblogs: a survey. *Data & Knowledge Engineering*, 53(3): 225-241, 2005.
- [21] M. Masseroli, D. Martucci, and F. Pinciroli. GFINDER: Genome Function INtegrated Discoverer through dynamic annotation, statistical analysis, and mining. *Nucleic Acids Research*, 32(Web Server issue):W293-W300, 2004.
- [22] M. Masseroli, O. Galati, M. Manzotti, K. Gibert, F. Pinciroli. Inherited disorder phenotypes: Controlled annotation and statistical analysis for knowledge mining from

- gene lists. *BMC Bioinformatics*, 6(Suppl 4):S18, 1-8, 2005.
- [23] M. Masseroli. Management and analysis of genomic functional and phenotypic controlled annotations to support biomedical investigation and practice. *IEEE Transaction on Information Technology in Biomedicine*, 11(4):376-385, 2007.
- [24] M. Masseroli, H. Kilicoglu, F-M. Lang, and T.C. Rindflesch. Argument-predicate distance as a filter for enhancing precision in extracting predications on the genetic etiology of disease. *BMC Bioinformatics*, 7(1):291, 1-12, 2006.
- [25] S. Ceri and P. Fraternali. Model for the definition of World Wide Web sites and methods for their design and evaluation. Patent US 6,591,271, July 2003.
- [26] S. Ceri, P. Fraternali, A. Bongio, M. Brambilla, S. Comai, and M. Matera. Designing data-intensive Web applications. Morgan-Kaufmann Series in Data Management Systems, J. Gray ed., Morgan-Kaufmann, 2003.
- [27] I. Manolescu, M. Brambilla, S. Ceri, S. Comai, and P. Fraternali. Model-driven design and deployment of service-enabled Web applications. *ACM Transactions on Internet Technology*, 5(3):439-479, 2005.
- [28] M. Brambilla, S. Ceri, S. Comai, and P. Fraternali. A CASE tool for modelling and automatically generating Web service-enabled applications. *Int. Journal of Web Engineering and Technology (IJWET)*, 2(4):354-372, 2006.
- [29] M. Brambilla, S. Ceri, P. Fraternali, R. Acerbis, and A. Bongio. Model-driven design of service-enabled Web applications. In ACM SIGMOD Conf. 2005, 851-856.
- [30] M. Brambilla, S. Ceri, P. Fraternali, and I. Manolescu. Process modelling in Web applications. *ACM Transactions on Software Engineering and Methodology*, 15(4):360-409, 2006.
- [31] M. Brambilla and C. Tziviskou. Fundamentals of exception handling within workflow-based Web applications. *Journal of Web Engineering*, 4(1):38-56, 2005.
- [32] M. Brambilla, S. Ceri, S. Comai, and C. Tzivisko. Exception handling in workflow-driven Web applications. In Proc. WWW 2005, May 10-14, 2005, Chiba, Japan, pp. 170-179.
- [33] M. Matera, A. Maurino, S. Ceri, and P. Fraternali. Model-driven design of collaborative Web applications. *Software-Practice & Experience*, 33:701-732, 2003.
- [34] S. Ceri, F. Daniel, and F. Facca. Modelling Web applications reacting to user behaviours. Special issue on Web dynamics. *Computer Networks*, 50(10):1533-1545, 2006.
- [35] S. Ceri, P. Dolog, M. Matera, and W. Nejdl. Adding client-side adaptation to the conceptual design of e-learning Web applications. *Journal of Web Engineering*, 4(1):21-37, 2005.
- [36] P. Fraternali, P.L. Lanzi, M. Matera, and A. Maurino. Model-driven Web usage analysis for the evaluation of Web application quality. *Journal of Web Engineering*, 3(2):124-152, 2004.
- [37] S. Ceri, M. Matera, F. Rizzo, and V. Demaldé. Designing data-intensive Web applications for content accessibility using Web marts. *Communications of the ACM*, 50(4):55-61, 2007.
- [38] M. Brambilla, I. Celino, S. Ceri, D. Cerizza, and E. Della Valle. Model-driven design and development of semantic Web service applications. *ACM Transactions on Internet Technology (TOIT)*, 8(1), 2008 (in press).
- [39] M. Brambilla, I. Celino, S. Ceri, D. Cerizza, E. Della Valle, and F.M. Facca. A software engineering approach to design and development of semantic Web service applications. In Proc. of the 5th Int. Semantic Web Conf. (ISWC 2006), Athens, GA, USA, November 5-9, 2006, LNCS 4273, 172-186.

Report on the First International Workshop on Ranking in Databases (DBRank'07)

Ihab F. Ilyas
University of Waterloo
Waterloo, Ontario, Canada
ilyas@cs.uwaterloo.ca

Gautam Das
University of Texas at Arlington
Arlington, Texas, USA
gdas@cse.uta.edu

ABSTRACT

This report summarizes the presentations, keynotes and discussions that took place during the first international workshop on ranking in databases (DBRank'07). The workshop was held on April 16, 2007, in conjunction with ICDE in Istanbul, Turkey.

1. INTRODUCTION

The International Workshop on Ranking in Databases (DBRank) focuses on the semantics, the modeling and the implementation of ranking and ordering in database systems and applications. In recent years, there has been a great deal of interest in developing effective techniques for ad-hoc search and retrieval in a variety of domains such as relational databases, document and multimedia databases, and scientific information systems. In particular, a large number of emerging applications require exploratory querying on such databases; examples include users wishing to search databases and catalogs of products such as homes, cars, cameras, restaurants, and photographs. Traditional database query languages such as SQL follow the Boolean retrieval model, i.e., tuples that exactly satisfy the selection conditions laid out in the query are returned—no more and no less. While extremely useful for the expert user, this retrieval model is inadequate for ad-hoc retrieval by exploratory users who cannot articulate the perfect query for their needs—either their queries are very specific, resulting in no (or too few) answers, or are very broad, resulting in too many answers.

To address the limitations of the Boolean retrieval model in these emerging ad-hoc search and retrieval applications, Top- k queries and ranking query results are gaining increasing importance. In fact, in many of these applications, ranking is an integral part of the semantics, e.g., keyword search, similarity search in multimedia as well as document databases. The increasing importance of ranking is directly derived from the explosion in the volume of data handled by current applications. The sheer amount of data makes it almost impossible to process queries in the traditional compute-

then-sort approach. Hence, ranking comes as a great tool for soliciting user preferences and data exploration.

Ranking imposes several challenges for almost all data-centric systems. In relational databases, large body of work has been recently proposed to support ranking as a first class construct through rank-aware algebra, ranking operators and new optimization frameworks that integrate ranking in plan enumeration and costing. There has been exciting recent work on automatic learning of appropriate ranking functions for database applications (e.g., based on adaptations of IR ranking functions to leverage dependency information in structured data), on designing expressive languages for user preferences modeling, on adaptation of keyword querying paradigms to relational databases, as well as on exciting new developments in new Top- k algorithms for relational, documents and multimedia databases. Ranking query results in semi-structured and XML databases has been also the focus of many recent contributions.

DBRank'07 solicited full and short papers that describe current research and work-in-progress efforts in enabling ranking in database systems.

2. LOGISTICS

The program committee of DBRank'07 consisted of 19 expert members from academia and industry. We had 23 submissions, each received at least two reviews. 5 full (8 pages) papers and 5 short (4 pages) papers were accepted by the program committee. The submission quality was very high and we wished we could squeeze more papers in the one-day program. However, we decided to accept only 10 papers to leave enough time for each paper to be properly presented and discussed. Each full paper was given a 30 minutes slot, while short papers were assigned 15 minutes slots.

In addition to paper presentations, we had keynotes by two distinguished scholars, Dr. Surajit Chaudhuri (Microsoft Research) and Prof. Gerhard Weikum (Max-Planck Institute for Informatics). For additional information, please refer to the URL of the workshop: <http://www.cs.uwaterloo.ca/dbrank2007>.

We were really happy with the discussions triggered by the keynotes and by the paper presentations. Based on a head-count, we had around 30 attendees throughout the day. In the following sections of this report, we briefly comment on the keynotes and the technical contributions presented in DBRank'07.

3. TECHNICAL CONTRIBUTIONS

Multiple aspects of ranking were discussed in 10 presentations. We roughly categorize the papers as follows:

Preference Specification Kenneth Ross from Columbia University presented two papers on preference specification, formalism and evaluation based on partial orders. The first short paper [5] identified several anomalies in the behavior of conventional notions of composition for preferences defined by strict partial orders. Kenneth showed how these anomalies can be avoided by defining a preorder that extends the given partial order, and by using the pair of orders to define order composition. The presentation included multiple examples to show the unintuitive results of composing strict partial orders using several composition methods, e.g., prioritized and Pareto composition. The second full paper [6] described a study of constraint formalisms for expressing user preferences as base facts in a partial order. The paper proposes a language that allows comparison and a limited form of arithmetic. The paper also shows that the transitive closure computation is required to complete the partial order terminates. Preference query processing was also briefly addressed in this paper, where index structures were presented to allow efficient evaluation over large data sets.

Scoring and Ranking Functions Scoring is a fundamental challenge in supporting ranking of database objects. While several rank aggregation and preference handling algorithms have been proposed, all these techniques depend on some sort of scores or a scoring function to be provided by the application or by the data generating process. Three papers focused on providing such scoring mechanism in different contexts.

Aparna Varde of Virginia State University presented a full paper [9] on learning the relative importance of similarity features in the context of image retrieval. The paper proposes a method called **FeaturesRank** for learning a distance function between the query image and images stored in the database. A training sample with pairs of images is used and the extent of similarity is identified for each pair. **FeaturesRank** clusters the given images in levels. It then adjusts the distance function based on the error between the clusters and training samples using multiple heuristics. **FeaturesRank** was evaluated with real image data from nanotechnology and bioinformatics.

Gultekin Ozsoyoglu of Case Western Reserve University presented a full paper [4] on comparing the quality

of scoring functions in the context of searching literature in digital libraries. The paper discusses three different functions that assign scores to papers based on their context. The extensive experimental study compares the quality of these functions based on accuracy (precision) and separability (uniformity of output scores), and shows that the text-based and the pattern-based scores yield better accuracy and separability than the citation-based scores.

Arthur Van Bunningen of the University of Twente presented a full paper [1] that proposes a novel explanatory approach of looking at context-aware relevance by defining the context-aware relevance of features as a probabilistic function of past choices. Context-aware preference is introduced as a way to capture the changes in user needs and preferences with respect to the search context. The paper shows that this approach goes well together with traditional probabilistic information retrieval and uncertainty of context information.

Ranking in XML Ranking is a natural way to explore and query large volumes of XML documents. In contrast to keyword search in text documents or ranking query results in relational database systems, the characteristics of XML data impose unique combination of ranking based on value and structural similarity. Two prototypes for enabling ranking in XML database were discussed in DBRank'07: the **ArHeX** and the **TReX** systems.

Ismael Sanz of the Universitat Jaume I, Spain presented a short paper [7] describing the features of the **ArHeX** similarity-oriented XML processing toolkit. **ArHeX** is designed to assist in the engineering of XML similarity-oriented applications, and to support the design and evaluation of suitable similarity measures and their associated indexes for each specific application.

Mariano Consens from the University of Toronto presented a full paper [2] that addresses retrieval queries that combine structural constraints with keyword search in XML database. The paper describes the **TReX** system an XML retrieval system that can exploit multiple structural summaries (including newly proposed ones). **TReX** can also self-manage small, redundant indexes to speed up the evaluation of workloads of top- k queries. The redundant indexes are maintained to enable **TReX** to select among different evaluation strategies.

Skyline Queries Skyline queries are natural ways to express ranking requirements in the absence of a concrete scoring function that aggregates the scores of multiple ranking criteria. Marcel Karnstedt of the Technische Universitt Ilmenau, Germany presented a short paper [3] that proposes three variants of a skyline operator and two extensions, especially suitable for efficient determination of skylines in structured overlays peer-to-peer environments.

Ranking in Other Domains Two short papers discussed the application of ranking in domains other than traditional relational and XML database systems. Jiawei Han of the University of Illinois at Urbana-Champaign presented a short paper [10] that addresses efficient evaluation of ranking queries in OLAP environments, introducing the *ranking cube*: a semi off-line materialization and semi-online computation model for answering top- k queries. On the other hand, Sharma Chakravarthy of the University of Texas at Arlington presented a short paper [8] that highlights multiple ranking issues in the context of information integration.

4. KEYNOTES

We were fortunate to have two outstanding keynote talks. Prof. Gerhard Weikum of Max-Planck Institute for Informatics reviewed the advantages and disadvantages of TA (the Threshold Algorithm) and its many extensions, putting them in perspective against algorithmic alternatives and pointing out unsolved technical issues and research opportunities. The family of Threshold Algorithms has become a very popular method for top- k query processing and ranked retrieval of unstructured, semi-structured, and structured data. TA has many elegant properties, such as instance optimality, and is extremely versatile. However, Gerhard's talk demonstrated that it also has specific limitations and is competing with alternative methods for top- k queries. The talk gave a nice overview and a timeline for the evolution of ranking algorithms in the last decade highlighting their strengths and weaknesses, and described recent work on optimizations and variations of TA in tackling several technical challenges. Examples of discussed challenges are choosing the best order for rank aggregation in nested top- k queries and relaxing ranking criteria.

The second keynote address was given by Dr. Surajit Chaudhuri of Microsoft Research who discussed various aspects of enabling general search over relational databases. The talk reviewed semantics and efficiency issues in supporting keyword search and ranking over databases, and critically examined past and current research. The talk highlighted important issues that require more attention in future research, e.g., (a) role of applications/business objects (b) architectural considerations - separation of functionality in database server vs. middleware. This interesting talk by Surajit intrigued the audience to engage in an active discussion debating the applicability of recent research contributions to real-world database engines. This unique industry perspective of Surajit raised many flags which will help shaping the emerging area of ranking in database systems.

5. FUTURE OF DBRANK

After receiving many unsolicited positive comments from the workshop attendees, we were encouraged to try making DBRank a yearly event. DBRank'08 will be held in conjunction with ICDE 2008 in Cancun, Mexico on the 11th and 12th of April. DBRank'08 is co-chaired by Vagelis Hristidis (Florida International University) and Ihab F. Ilyas (University of Waterloo). We are looking forward to having another interesting and successful round of DBRank.

6. ACKNOWLEDGMENTS

We would like to express our sincere gratitude to our distinguished program committee members who provided timely, high quality and constructive reviews to all the papers. We were really impressed by their enthusiasm and dedication to make DBRank'07 a great success. This workshop would not have been successful without the huge support by our two keynote speakers who gave exciting talks full of research challenges and future research directions to this emerging community. Last but not least, we would like to thank all the authors who submitted their self-selected high quality papers to DBRank'07.

7. PAPERS PRESENTED IN DBRANK'07

- [1] Arthur Van Bunningen, Maarten Fokkinga, Peter Apers, and Ling Feng. Ranking query results using context-aware preferences.
- [2] Mariano Consens, Xin Gu, Yaron Kanza, and Flavio Rizzolo. Self managing top-k (summary, keyword) indexes in xml retrieval.
- [3] Marcel Karnstedt, Jessica Muller, and Kai-Uwe Sattler. Cost-aware skyline queries in structured overlays.
- [4] Nattakarn Ratprasartporn, Sulieman Bani-Ahmad, Ali Cakmak, Jonathan Po, and Gultekin Ozsoyoglu. Evaluating different ranking functions for context-based literature search.
- [5] Kenneth Ross. On the adequacy of partial orders for preference composition.
- [6] Kenneth Ross, Peter Stuckey, and Amelie Marian. Practical preference relations for large data sets.
- [7] Ismael Sanz, Rafael Berlanga, Marco Mesiti, and Giovanna Guerrini. Flexible composition of indexes and similarity measures in xml.
- [8] Aditya Telang, Roochi Mishra, and Sharma Chakravarthy. Ranking issues for information integration.
- [9] Aparna Varde, Elke Rundensteiner, Giti Javidi, Ehsan Sheybani, and Jianyu Liang. Learning the relative importance of features in image data.
- [10] Dong Xin and Jiawei Han. Integrating olap and ranking: The ranking-cube methodology.

Report on the Fourth International Workshop on Data Management for Sensor Networks (DMSN 2007)

Magdalena Balazinska
Univ. of Washington
magda@cs.washington.edu

Amol Deshpande
Univ. of Maryland
amol@cs.umd.edu

Alexandros Labrinidis
Univ. of Pittsburgh
labrinid@cs.pitt.edu

Qiong Luo
Hong Kong Univ. of Sci. & Tech.
luo@cse.ust.hk

Samuel Madden
MIT
madden@csail.mit.edu

Jun Yang
Duke Univ.
junyang@cs.duke.edu

Introduction

Sensor networks enable an unprecedented level of access to the physical world, and hold tremendous potential to revolutionize many application domains. Research on sensor networks spans many areas of computer science, and there are now major conferences, e.g., *IPSN* and *SenSys*, devoted to sensor networks. However, there is no focused forum for discussion of early and innovative work on data management in sensor networks. The *International Workshop on Data Management for Sensor Networks (DMSN)*, inaugurated in 2004, aims to fill this significant gap in the database and sensor network communities.

DMSN 2007, the fourth DMSN workshop, was held on September 24, 2007 in conjunction with the *33rd International Conference on Very Large Database (VLDB 2004)* in Vienna, Austria. Building on the successes of the three previous DMSN workshops, DMSN 2007 aimed at bringing together researchers working on all aspects of sensor data management: from data processing in networks of remote, wireless, resource-constrained sensors to managing heterogeneous, noisy, and sometimes sensitive sensor data in databases. The resource-constrained, lossy, noisy, distributed, and remote nature of sensor networks means that traditional database techniques often cannot be applied without significant re-tooling. Challenges associated with acquiring and processing large-scale, heterogeneous sets of live sensor data also call for novel data management techniques. Finally, in many applications, collecting sensor data raises important privacy and security concerns that require new protection and anonymization techniques.

As the field of sensor networks continues to develop, we have expanded the scope of DMSN 2007 from previous workshops in the series, by encouraging contributions on a broader set of topics, including: database languages

for sensor tasking; distributed sensor data storage and indexing; data replication and consistency in noisy and lossy environments; energy-efficient data acquisition and dissemination; in-network query processing; networking support for data processing; query optimization and deployment planning in sensor networks; database techniques for managing loss, uncertainty, noise, and ambiguity; model-based sensor data processing; challenges and techniques for new types of sensor data, e.g., RFID, images and videos, data from scientific and medical instruments; personal, ubiquitous applications of sensor-based infrastructures; integration of sensor data of different modalities and from different sources; integration of sensor data in traditional databases and streaming systems; techniques for secure sensor data collection and processing; and privacy protection techniques for sensor data.

Program

The workshop program includes a keynote address, three research sessions with a total of seven papers, and a panel discussion on the present and future of sensor data management research. As a response to the Call for Papers this year, we received 15 full paper submissions. During the review process, each paper was reviewed by three or four members of the program committee or external reviewers and was also carefully discussed, resulting in the acceptance of seven papers.

Keynote Address The keynote address was delivered by Prof. Gustavo Alonso from ETH Zurich, with the title: “Myths and Realities of Sensor Network Data Management.” After about a decade of research on sensor networks, there is a growing body of expertise and experience in real deployments, and a number of companies have begun to commercialize sensor networks. It is now the high time that we reexamine our research efforts in perspective. Gustavo started his talk by questioning the assumptions about wireless sensor

networks that are commonly made by the research community, and gave the audience a healthy dose of reality check: Over and over again, real deployments show that some of the claimed properties of sensor networks (easy deployment, low cost, lack of alternative solutions, self organization, large-scale monitoring) are still not attainable in fields today. Gustavo then proposed a number of more realistic, but nevertheless ambitious, targets for research in sensor networks. First, existing programming tools for sensor networks still require considerable expertise and provide little support for reliability; for sensor networks to make a real impact, we must provide tools for real users instead of novelties for nerds, and develop turnkey solutions that address all aspects of an application, including production, deployment, maintenance, and reuse, as well as data collection, storage, cleansing, and analysis. Gustavo also argued for an approach to system design that is driven by real applications and deployments. There exists no one-size-fits-all solution; instead, we need to evaluate alternative designs and techniques in the context of each application scenario's requirements and constraints. Gustavo's call for tackling real, non-tutored, long-term deployments of sensor networks and linking them to the rest of the IT infrastructure was met with much support from audience.

Research Sessions The seven accepted papers are organized into three sessions: 1) in-network processing, 2) novel sensing modalities, and 3) modeling and programming sensor networks.

The first session features three papers on in-network processing, an idea that has been proven effective in reducing the cost of communicating data to base stations. The paper presented by Shili Xiang (National Univ. of Singapore) shows how to optimize a large number of queries across multiple base stations attached to the same sensor network, exploiting similarities among queries allocated to the same base station during in-network processing. The paper by Yongxuan Lai (Renmin Univ. of China) et al. proposes an in-network data-centric storage scheme that is dynamically balanced to avoid hotspots caused by data skewness, as well as to reduce the cost of data reporting and storage. Finally, the paper presented by Demetrios Zeinalipour-Yazti (Open Univ. of Cyprus) investigates how a swarm of moving sensor nodes can collaboratively compute a perimeter; data is acquired from the perimeter and aggregated and replicated by nodes inside the perimeter for reliability.

The second session features two papers that investigate novel sensing modalities that are much more complex than numeric readings such as temperature and light. Josh Hyman (UCLA) presented work on modeling the

rate of CO₂ uptake of drought-tolerant moss using a series of images taken of the plant; the work shows that regression tree models based on color features of the images taken by inexpensive image sensors can effectively predict CO₂ uptake, which is very difficult to measure directly. The second paper, presented by Saket Navlakha (Univ. of Maryland), studies vehicle tracking using video streams produced by traffic cameras; their novel approach, which represents image feature as a graph capturing neighboring relationships among vehicles, makes tracking feasible even with low-quality video streams.

In the third session on modeling and programming sensor networks, Pablo Guerrero (TU Darmstadt) made a case for using workflows to express the application logic of a wireless sensor and actuator network. This approach is better at capturing the actuation aspect than previously proposed programming paradigms, which focused primarily on sensor data collection. The second paper, presented by Yanif Ahmad (Brown Univ.), proposes a framework called *Pulse* for processing continuous queries over continuous-time data models, appropriate for many physical phenomena monitored by sensor networks; optimizations possible under *Pulse* can lead to significant performance savings compared with the traditional strategy of processing sensor data as streams of discrete tuples.

Panel Discussion The workshop concluded with a roundtable discussion featuring four panelists: Minos Garofalakis (Yahoo! Research and UC Berkeley), Zachary Ives (Univ. of Pennsylvania), Samuel Madden (MIT), and Sunil Prabhakar (Purdue Univ.). All panelists answered a resounding “no” to the question posed in title of the panel: “Sensor Data Management: Are We Done?”

Sam expressed optimism on the impact of database research on sensor network applications, but suggested that the DMSN community should reach out further—not only to the larger sensor networking and the more general distributed systems communities, but also to domain-specific communities with sensor applications. Echoing Gustavo's point in the keynote address, Sam cautioned against unrealistically simple simulations as well as overly ambitious deployment projects based on unproven technology.

The database community today is still mostly wedded to exact answers. Minos, Sunil, and Zack all identified support for approximation, uncertainty, and statistical methods as an important challenge. The inherently noisy and incomplete nature of sensor data makes sensor data management an ideal application domain for emerging probabilistic and uncertain databases.

Minos also highlighted distributed stream processing as an area with many open challenges motivated by sensor data management, including distribution of data and processing, and support for complex “queries” such as probabilistic inference.

Zack likened in-network sensor data processing to the old problem of distributed query processing, but with “extreme” degrees of distribution and constraints that make the problem more interesting. He also pointed out an analogous connection between sensor data fusion and the data integration problem.

Past and Future

The series of DMSN workshops has proved very successful at its goal of providing a forum for database researchers to present their innovative ideas on sensor data management. Over the last four years, we have had a collection of great papers on topics ranging from in-network storage and querying to building middleware systems for processing sensor data. We have also been fortunate to have had excellent keynotes and invited talks - Dr. Wei Hong from Intel Research (2004), Dr. Henry Tirri from Nokia Research (2005), Prof. Mike Franklin from UC Berkeley (2006), and Prof. Gustavo Alonso from ETH Zurich (2007). The themes of the workshop over the years mirror the development of sensor data management field itself. In 2004 and 2005, the workshops focused more on understanding the needs of sensor network applications and the role of sensor data management, as evidenced by the talks by Wei Hong and Henry Tirri. In the last two years, however, the focus shifted to higher-level issues such as data streams, uncertain data management, and infrastructures for large-scale data management and sharing.

Research in sensor data management is at a crossroads. There has been much progress in solving the problems, but many hard research challenges remain open, as highlighted by the participants of this year’s panel. We see an ongoing need for a forum like DMSN. The interest in the workshop remains high, as witnessed by the consistently high registration numbers and paper submissions from all over the world. Hence, we plan to continue to organize DMSN workshops in the near future.

Acknowledgements

Many people have contributed to the success of DMSN 2007. We would like to thank all authors who have submitted papers, as well as members of the program committee and the external reviewers; without them, it would have been impossible to put together a high-

quality program for DMSN 2007. We would like to thank the Conference Management Toolkit (CMT) Support Team at Microsoft Research for making CMT available to us and Intel Corporation for their generous support of the workshop. Finally, we would like to thank the local conference organization chair, Bernhard Schandl, and the VLDB 2007 workshop organizing committee—in particular Georg Gottlob, H. V. Jagadish, and Krithi Ramamritham—for their help.

Report on the first VLDB workshop on Management of Uncertain Data (MUD)

Ander de Keijzer[•] Maurice van Keulen[•]
[•]University of Twente
The Netherlands
{a.dekeijzer,m.vankeulen}@utwente.nl

Alex Dekhtyar[◊]
[◊]California Polytechnic State University
United States
dekhtyar@csc.calpoly.edu

1 Introduction

On Monday September 24th, we organized the first international VLDB workshop on Management of Uncertain Data [dKvKD07]. The idea of this workshop arose a year earlier at the Twente Data Management Workshop on Uncertainty in Databases [dKvK06]. The TDM is a bi-annual workshop organized by the Database group of the University of Twente, for which each time a different topic is chosen. The participants of TDM 2006 were enthusiastic about the topic “Uncertainty in Databases” and strongly expressed the wish for a follow-up co-located with an international conference. To fulfill this wish, we organized the MUD-workshop at VLDB.

The program committee consisted of 20 members and 1 advisory member, Jennifer Widom from Stanford University. Committee members came from universities and research institutes from Europe and North America. We accepted 6 full papers and invited 2 speakers, Lise Getoor from the University of Maryland at College Park and Sunil Prabhakar from Purdue University.

Both the morning and afternoon session consisted of an invited talk and a research session. In the morning Lise Getoor gave a talk on *Combining Tuple and Attribute Uncertainty in Probabilistic Databases*, which was followed by a research session on *Applications of Uncertain Data*. The afternoon session started with the talk by Sunil Prabhakar on *Supporting Probabilistic Data in Relational Databases*, which was followed by a session on *Querying Uncertain Data*.

Special thanks go to the Centre for Telematics and Information Technology (CTIT) for sponsoring the proceedings.

2 Applications of Uncertain Data

The kick-off of the workshop was given by Lise Getoor from the University of Maryland. She gave an overview of techniques from machine learning and reasoning under uncertainty. These areas have devel-

oped quite powerful models for the representation of probability distributions, for example, probabilistic graph models. She showed how these techniques influenced her work on probabilistic relational models especially on how to unify attribute and record level uncertainty.

The first research talk of the workshop was given by Antoon Bronselaer from the University of Ghent. He introduced the application of disaster victim identification for large scale disasters. The problem can be seen as an object matching or entity resolution problem: based on the available data of a victim determine whether or not that data refers to the same real world object as data from a reference list. The focus of the paper was on how to integrate a complex matching technique based on ear biometrics into their object matching framework. It was shown that the framework, which was based on a possibilistic uncertainty model, was capable of effectively capturing and handling the uncertainty resulting from missing data and feature extraction errors.

In the second presentation, Matteo Magnani argued that data integration could be the killer application for uncertain data management systems. One of the main problems in data integration is schema matching. Current approaches combine the judgments of multiple matchers to obtain the most relevant schema mappings. Magnani argues that significant improvement can be obtained by not only finding the correct mappings, but also by managing the incorrect ones properly. They propose to view the mappings as possible mappings with a certain level of uncertainty and treating the accompanying data during querying accordingly, i.e., also with a certain level of uncertainty.

The topic of the third and last presentation of the morning session was fuzzy querying. Ramón Alberto Carrasco presented their language dmFSQL (data mining fuzzy SQL) which allows you to easily verify data mining hypotheses. The paper focused on their latest addition to the language: fuzzy global dependencies. The idea is that the system computes the percentage of tuples which fulfill a given antecedent and consequent together w.r.t. those that only fulfill the consequent. This allows you to validate hypoth-

esized monotonicity of relationships between objects in the data, e.g., which patterns imply higher earnings of a specific share on the stock market. For this particular example, it was presented how the system could obtain the final statement “Greater williams index and roughly equal moving average implies a greater value for the specific enterprise Telefonica with confidence 0.9”.

3 Querying Uncertain Data

The afternoon session started with an invited talk by Sunil Prabhakar. The topic of the talk was Supporting Probabilistic Data in Relational Databases and was focused on the ORION DBMS. Sunil Prabhakar provided a nice overview of possible world semantics and the problems that arise with continuous uncertainty. Currently, the ORION system offers the combination of continuous uncertainty and possible world semantics.

The first research talk of the afternoon session on Querying Uncertain Data was given by Patrick Bosc from IRISA/ENSSAT, France. The model of uncertain data he used was a possibilistic model. During the discussion at the end of the talk, many of the questions were addressing the differences between probabilistic and possibilistic theory. One notable difference is that a possibilistic model uses maximum and minimum for combining confidences, while a probabilistic model uses addition and multiplication. From the discussion arose that a possibilistic model does not make assumptions about dependencies between stochastic variables while probabilistic models usually so. The conclusion of the discussion was, that both theories have their advantages and purposes.

The second presentation, given by Jef Wijsen of the University of Mons-Hainaut, was about introducing uncertainty by considering possible repairs for key constraint violations. These violations can be solved in different ways. Each of the minimal solutions can be regarded as a possible world. Jef focused on the notion of relations to ‘consistently join’, i.e., for all possible repairs the join contains at least one tuple. He used a game theoretic approach to decide on this notion.

The last presentation of the workshop was given by Raghobham Murthy of Stanford University. His presentation on aggregate functions in databases supporting uncertainty used the Trio database system as an example. He presented algorithms for estimating a lower bound, higher bound, and expected value for aggregates on uncertain relations, because these typically produce exponential results. Afterwards, even after the workshop officially ended, several participants continued discussing about the semantics of aggregates. Different views on how aggregates should be interpreted were discussed, and

in the end it turned out that people agreed on the main idea of aggregates, although there seemed to be some difference in opinions on the details. All in all, this topic will probably be continued at subsequent workshops.

4 Conclusion and Outlook

Discussions during the workshop showed that the management of uncertain data is a vibrant research area with many promising applications, but also a significant number of open issues. For example, one can distinguish several kinds of underlying data models for uncertain data: fuzzy logic-based models, repair models, possibilistic models and probabilistic models. The relationships, commonalities and differences are not well understood yet. And if theory is not well enough established yet, work on algorithms, scalability and systems is necessarily also still in its infancy. But, the strength of approaches based on properly managing uncertainty in data can already be demonstrated as the application-oriented papers in the MUD workshop clearly show. Moreover, the papers in this workshop also show that the challenges, for example the ones presented in the visionary paper on dataspace systems [HFM06], are being addressed today and significant advances are being made. To continue our efforts to build a rich co-operating community on this topic and support effective exchange of ideas, we plan to organize a second MUD-workshop again co-located with VLDB next year.

References

- [dKvK06] A. de Keijzer and M. van Keulen, editors. *Proc. of the 2nd Twente Data Management Workshop (TDM 2006) on Uncertainty in Databases (Enschede, The Netherlands, June 6, 2006)*, number WP06-01 in CTIT Workshop Proceedings Series. Centre for Telematics and Information Technology (CTIT), Univ. of Twente, Enschede, The Netherlands, June 2006. <http://www.cs.utwente.nl/~tdm>.
- [dKvKD07] A. de Keijzer, M. van Keulen, and A. Dekhtyar, editors. *Proc. of the 1st Int. VLDB workshop on Management of Uncertain Data (MUD, Vienna, Austria, September 24, 2007)*, number WP-CTIT-07-08 in CTIT Workshop Proceedings Series. Centre for Telematics and Information Technology (CTIT), Univ. of Twente, Enschede, The Netherlands, September 2007. <http://mud.cs.utwente.nl>.
- [HFM06] A.Y. Halevy, M.J. Franklin, and D. Maier. Principles of dataspace systems. In *Proceedings of PODS, Chicago, IL, USA*, pages 1–9, 2006.

2008 ACM SIGMOD AWARDS

The SIGMOD Awards Committee is now open to receiving nominations for the SIGMOD Edgar F. Codd Innovations Award and the SIGMOD Contributions Award. In addition, the awards committee welcomes informal advice regarding the Test of Time Award for the paper in the 1998 SIGMOD Conference that has had the most impact since its publication.

The nomination deadline is **April 7, 2008**.

In 1992, ACM SIGMOD started the Annual SIGMOD Innovations Award and SIGMOD Contributions Award as part of its Awards Program. In 2004, SIGMOD, with the unanimous approval of ACM Council, renamed the Innovations Award in honor of Dr. Edgar F. (Ted) Codd (1923 - 2003), who invented the relational data model and was responsible for the significant development of the database field as a scientific discipline. The previous winners of the Innovations Award are: Michael Stonebraker (1992), James Gray (1993), Philip Bernstein (1994), David DeWitt (1995), C. Mohan (1996), David Maier (1997), Serge Abiteboul (1998), Hector Garcia-Molina (1999), Rakesh Agrawal (2000), Rudolf Bayer (2001), Patricia Selinger (2002), Donald Chamberlin (2003), Ronald Fagin (2004), Michael Carey (2005), Jeffrey Ullman (2006), and Jennifer Widom (2007). The previous winners of the Contributions Award are: Maria Zemankova (1992), Gio Wiederhold (1993), Yahiko Kambayashi (1995), Jeffrey Ullman (1996), Avi Silberschatz (1997), Won Kim (1998), Raghu Ramakrishnan (1999), Laura Haas and Michael Carey (2000), Daniel Rosenkrantz (2001), Richard Snodgrass (2002), Michael Ley (2003), Surajit Chaudhuri (2004), Hongjun Lu (2005), Tamer Ozsu (2006), and Hans-Joerg Schek (2007).

INNOVATIONS/CONTRIBUTIONS AWARDS

- | | | |
|--------------|--|--|
| 1. Name: | SIGMOD Edgar F. Codd Innovations Award | SIGMOD Contributions Award |
| 2. For What: | Innovative and highly significant contributions of enduring value to the development, understanding, or use of database systems and databases. | Outstanding and sustained services to the database field through education, conference organizations, journals, standards activities, research funding, etc. |
3. Given: Annually (if there is at least one qualified candidate).
 4. Award: A plaque per person plus \$1000 per award (the latter to be split among a group, if it is a group award).
 5. Administration: Administered by the SIGMOD Awards Committee.
 6. Nomination/Evaluation Procedures: Anyone in the field can nominate one or more persons or groups (self nominations are excluded). Nominations should include a proposed citation (up to 25 words), a succinct (100-250 words) description of the innovation/contribution, and a detailed statement to justify the nomination; plaintext is preferred. At least three additional supporting letters should be submitted. Such letters, however, should not be simple endorsements of the nomination, but convey additional factual information. The Awards Committee will evaluate all nominations and decide on zero or more winners. Nominations must be received by April 7, 2008 to be considered for this year's award.
 7. Recipients: The recipients will receive the awards at the annual ACM SIGMOD/PODS Conference, at the awards luncheon; each awardee will give a short speech (5-10 minutes).
 8. Eligibility for Nomination: Anyone except the current elected officers of SIGMOD (Chair, Vice Chair, and Treasurer), and members of the SIGMOD Awards Committee. Awards should be for contributions not already honored by a major ACM Award (e.g., the Turing Award, SIGMOD Edgar F. Codd Innovations Award, or SIGMOD Contributions Award).

TEST-OF-TIME AWARD

There is no formal nomination process for this award, but input from the database research community is welcome. The SIGMOD Awards Committee is charged with selecting the paper from the SIGMOD Proceedings from 10 years ago (i.e., from 1998 SIGMOD, for this year) that has best met the "test of time," that is, it has had the most influence since its publication. We are especially interested in first-hand accounts of ways in which the ideas of a paper have been used in practice. Take a look at the 1998 SIGMOD Proceedings (<http://www.sigmod.org/sigmod/dblp/db/conf/sigmod/sigmod98.html>), and if you have any information you believe would be of use to the committee, then please send the committee a note, as described below.

WHERE TO SEND NOMINATIONS

Nominations should be submitted via e-mail to the chair of the SIGMOD Awards Committee, Gerhard Weikum, with copies to the other members of the committee.

SIGMOD Awards Committee:

- Peter Buneman
University of Edinburgh
`opb@inf.ed.ac.uk`
- Michael J. Carey
BEA Systems, Inc.
`mcarey@bea.com`
- David Maier
Portland State University
`maier@cs.pdx.edu`
- Laura Haas
IBM Almaden Research Center
`laura@almaden.ibm.com`
- Gerhard Weikum (Chair)
Max-Planck Institute of Computer Science, Saarbruecken
`weikum@mpi-sb.mpg.de`

Call for Submissions

ACM SIGMOD 2008 Undergraduate Research Poster Competition

SIGMOD/PODS 2008 Conference – June 9-12, 2008, Vancouver, Canada

Chair: Lukasz Golab
AT&T Labs-Research
lgolab@research.att.com

This year's SIGMOD conference will give undergraduate students an opportunity to showcase their research accomplishments in a poster competition. Up to five students will be selected to attend the conference and present posters to other attendees of SIGMOD/PODS 2008. For each invited student, up to US\$1000 will be provided to defray conference attendance costs (registration fee, travel, lodging, etc). A "best poster" winner will be selected by the competition chair and announced at the SIGMOD 2008 awards session.

Undergraduate students who have played a key role in a research project are invited to submit an abstract to the poster competition. Any research projects broadly related to data management are within the scope of the competition (for a list of sample areas of interest, see the SIGMOD call for papers at: http://www.sigmod08.org/sigmod_call_papers.shtml).

Based on the abstracts, the competition chair will choose up to five students to invite to the SIGMOD/PODS conference and present posters. For the purposes of this competition, a student is considered an undergraduate student if he/she has not yet obtained a BS (or equivalent) degree or has obtained that degree on or after December 2007, and he/she is not enrolled in a graduate program at the time of submission. If the applicant's school system is "non-traditional", and the applicant considers him/herself eligible, then the competition chair should be contacted before an abstract is submitted.

Submission Guidelines

In order to submit an abstract to the research poster competition, students must send an email to the competition chair (lgolab@research.att.com) by **Friday, April 4, 2008, 5pm PST**. The subject of the email must be "<candidate's full name> SIGMOD UNDERGRADUATE POSTER COMPETITION". The following information must be included (not attached) in the email in plain text. No HTML, PDF, Postscript or any other formats will be accepted.

1. Name of department and school, and current academic status, including the number of years until graduation.
2. Name of academic advisor.
3. An abstract of up to 800 words explaining the proposed content of the poster, including:
 - (a) a clear and concise problem statement,
 - (b) brief technical overview of the solution,
 - (c) summary of major results (e.g., "faster than existing solutions by x percent").
4. Description of the role played by the student in the project.

All submissions must be in plain text with the proper subject line as explained above. Any submission that does not satisfy these conditions may be flagged as junk mail and automatically discarded without further notification. Decisions will be emailed by Monday, April 14, 2008; authors of accepted abstracts will receive further instructions at that time. The competition chair reserves the right to reject all submissions.

Note: submissions to the research poster competition are permitted even if the student already has a paper on the same topic due to appear at the SIGMOD/PODS 2008 conference.

Important Dates

- Submission deadline: Friday, April 4, 2008, 5pm PST
- Notification of results: Monday, April 14, 2008

Comments and questions should be directed to the competition chair at lgolab@research.att.com

Tribute to Honor Jim Gray

May 31, 2008 @ UC Berkeley

<http://www.eecs.berkeley.edu/IPRO/JimGrayTribute>

Jim Gray's family, friends and colleagues have arranged a public day of talks and reminiscences in his honor, to be held on May 31, 2008, at UC Berkeley. All are welcome.

Although Dr. Gray will be officially listed as missing until 2012, his family has asked that we have this Tribute now, to honor him before too much time has passed.

There are two parts to the Tribute. The General Session that begins the day is intended for the general public. The Technical Sessions will go through the afternoon, and are intended for a computer science audience. The organizers request that people wishing to attend the Technical Sessions register in advance on the web, to facilitate planning.

More information, including the program for the day, registration forms, and information on accommodations is available on the web at <http://www.eecs.berkeley.edu/IPRO/JimGrayTribute>. Questions about event logistics may be directed to jimgraytribute@eecs.berkeley.edu.

The search for Dr. Gray is detailed in Wired Magazine's article:
http://www.wired.com/techbiz/people/magazine/15-08/ff_jimgray

Call For Papers

Fourth International Workshop on Data Management on New Hardware (DaMoN 2008)

Colocated with
ACM SIGMOD/PODS 2008

Vancouver, Canada
June 13, 2008



Objective

The aim of this one-day workshop is to bring together researchers who are interested in optimizing database performance on modern computing infrastructure by designing new data management techniques and tools.

Motivation

The continued evolution of computing hardware and infrastructure imposes new challenges and bottlenecks to program performance. As a result, traditional database architectures that focus solely on I/O optimization increasingly fail to utilize hardware resources efficiently. CPUs with superscalar out-of-order execution, simultaneous multi-threading, multi-level memory hierarchies, and future storage hardware (such as MEMS) impose a great challenge to optimizing database performance. Consequently, exploiting the characteristics of modern hardware has become an important topic of database systems research.

The goal is to make database systems adapt automatically to the sophisticated hardware characteristics, thus maximizing performance transparently to applications. To achieve this goal, the data management community needs interdisciplinary collaboration with computer architecture, compiler and operating systems researchers. This involves rethinking traditional data structures, query processing algorithms, and database software architectures to adapt to the advances in the underlying hardware infrastructure.

Important Dates (*tentative*)

Paper submission: **April 11**
Notification of acceptance: **May 2**
Camera-ready copies due: **May 16**

Topics Of Interest

We seek submissions bridging the area of database systems to computer architecture, compilers, and operating systems. In particular, submissions covering topics from the following non-exclusive list are encouraged:

- database algorithms and data structures on modern hardware
- cost models and query optimization for novel hierarchical memory systems
- hardware systems for query processing
- data management using co-processors
- query processing using computing power in storage systems
- database architectures for low-power computing and embedded devices
- database architectures on multi-threaded and chip multiprocessors
- performance analysis of database workloads on modern hardware
- compiler and operating systems advances to improve database performance
- new benchmarks for microarchitectural evaluation of database workloads

Organization

Workshop Co-Chairs

Kenneth Ross, Columbia University
Qiong Luo, HKUST

Program Committee

Anastasia Ailamaki, Carnegie Mellon University
Bishwaranjan Bhattacharjee, IBM Research
Peter Boncz, CWI Amsterdam
Shimin Chen, Intel Research
Goetz Graefe, HP Labs
Stavros Harizopoulos, HP Labs
Martin Kersten, CWI Amsterdam
Bongki Moon, University of Arizona
Jun Rao, IBM Research
Jingren Zhou, Microsoft Research



<http://www.cse.ust.hk/damon2008>

Call for Papers **DBTest 2008**

1st International Workshop on Testing Database Systems

Vancouver, June 13, 2008

<http://research.microsoft.com/dmx/dbtest2008>

Co-located with ACM SIGMOD Conf.

Motivation and Scope

The functionality provided by modern database management systems (DBMS) is continuously expanding. New trends in hardware architectures, new data storage requirements, and new usage patterns drive the need for continuous innovation and expansion in modern database engines. As a result, DBMS are becoming increasingly complex and difficult to validate. Moreover, while DBMS functionality has advanced significantly during the past 10 years, the methodology for testing and tuning has not evolved accordingly. Testing and tuning a database system are becoming increasingly expensive and are often dominating the release cycle of a database product. It is not unusual that fifty percent of the development cost is spent on testing and tuning and that several months are reserved for testing before a new release can be shipped. Without revolutionary new ideas, the situation is going to become even worse in future.

The purpose of this workshop is to expose to the academic community the challenges and practical impact of adequate database testing, and encourage further research in the area. The long term goal is to devise new technique that reduce the cost and time to test and tune database products so that users and vendors can spend more time and energy on actual innovations. Obviously, the general area of testing has attracted a great deal of attention in the software engineering community. However, testing DBMS imposes particular challenges and opportunities which have not been addressed in either the database or software engineering community. Only recently, DBMS testing has gained more attention in the database community.

The participants of this workshop will be from both industry and academia. In addition to novel techniques, the workshop will present war stories in order to define and better understand the problem space.

Topics of Interest

- DBMS testing techniques
- Generation of synthetic data for test databases
- Generation of stochastic test models for large test matrices

- Techniques and algorithms for automatic program verification
- Maximizing code coverage of engine components
- Testing correctness of DBMS components
- Test-modeling of DBMS engines and components
- Testing and designing systems that are robust to estimation inaccuracies
- Testing the efficiency of adaptive policies and components
- Minimizing, automating and ranking of engine tuning parameters
- Identifying performance bottlenecks
- Workload characterization with respect to performance metrics
- Workload characterization with respect to engine components
- Metrics for predictability of query and workload performance
- Metrics for query plan robustness
- Security and vulnerability testing
- War Stories

Paper Submission

DBTest 2008 invites the submission of original contributions in the area of database system testing and tuning. As mentioned above, DBTest is also interested in war stories and practitioners' reports on techniques and issues in testing and tuning database systems. Papers should be formatted according to the ACM guidelines and SIGMOD proceedings template available at:

http://www.sigmod08.org/sigmod_formatting.shtml

Papers should not be longer than six pages and should be submitted in PDF by E-Mail to the workshop chairs:

Leo Giakoumakis, Microsoft Corporation, USA (leogia@microsoft.com)

Donald Kossmann, ETH Zürich, Switzerland (kossmann@inf.ethz.ch)

Important Dates

Paper Submission:	April 11, 2008
Notification of acceptance:	May 9, 2008
Camera-ready:	May 23, 2008
Workshop:	June 13, 2008

Contact

Leo Giakoumakis, Microsoft Corporation, USA (leogia@microsoft.com)

Donald Kossmann, ETH Zürich, Switzerland (kossmann@inf.ethz.ch)



CALL FOR PAPERS

MobiDE 2008: Seventh International ACM Workshop on Data Engineering for Wireless and Mobile Access

June 13, 2007, Vancouver, Canada (in
conjunction with SIGMOD/PODS 2008)

<http://www.cs.ucy.ac.cy/mobide08/>



In-cooperation with



Important Dates:

March 19: Abstracts
March 26: Papers
April 30: Notification
May 14: Camera Ready

(all deadlines are midnight EST)

General Chairs:

Alex Delis
University of Athens
ad@di.uoa.gr

Vladimir I. Zadorozhny
University of Pittsburgh
vladimir@sis.pitt.edu

Program Chairs:

Yannis Kotidis
Athens University of
Economics and Business
kotidis@aueb.gr

Pedro Jose Marron
University of Bonn
pjmarron@cs.uni-bonn.de

Publicity Chair:

Demetris Zeinalipour
Open University of Cyprus
zeinalipour@ouc.ac.cy

MobiDE'08 is the seventh of a successful series of workshops that aims to act as a bridge between the data management, wireless networking, and mobile computing communities.

The 1st MobiDE workshop took place in Seattle (August 1999), in conjunction with MobiCom 1999; the 2nd MobiDE workshop took place in Santa Barbara (May 2001), together with SIGMOD 2001; the 3rd MobiDE workshop took place in San Diego (September 2003), together with MobiCom 2003; the 4th MobiDE workshop was held in Baltimore (June 2005). In 2006, MobiDE was organized in Chicago, IL (June 2006) and last year it was held in Beijing, China (June 2007). Since 2005, the event has been collocated with the annual SIGMOD conference. This year, MobiDE 2008 is sponsored by ACM SIGMOD and held in co-operation with ACM SIGMOBILE (pending approval).

The workshop will serve as a forum for researchers and technologists to discuss the state-of-the-art, present their contributions, and set future directions in data management for mobile and wireless access.

The topics of interest related to mobile and wireless data engineering include, but are not limited to:

- * ad-hoc networked databases
- * consistency maintenance and management
- * context-aware data access and query processing
- * data caching, replication and view materialization
- * data publication modes: push, broadcast, and multicast
- * data server models and architectures
- * database issues for moving objects: storing, indexing, etc.
- * m-commerce
- * mobile agent models and languages
- * mobility-aware data mining and warehousing
- * mobile database security
- * mobile databases in scientific, medical and engineering applications
- * mobile peer-to-peer applications and services
- * mobile transaction models and management
- * mobile web services
- * mobility awareness and adaptability
- * pervasive computing
- * prototype design of mobile databases
- * quality of service for mobile databases
- * sensor network databases
- * transaction migration, recovery and commit processing
- * wireless multimedia systems
- * wireless web



WebDB 2008



ACM SIGMOD/PODS 2008



11th International Workshop on the Web and Databases (WebDB 2008)

Friday, June 13, 2008 / Vancouver, Canada
(co-located with ACM SIGMOD/PODS 2008)
<http://webdb2008.com.polimi.it>

Call for Papers (Draft)

AIMS & TOPICS OF INTEREST:

The WebDB workshop focuses on providing a forum where researchers, theoreticians, and practitioners can share their knowledge and opinions about problems and solutions at the intersection of data management and the Web.

This year WebDB will focus on Web2.0 and Multimedia support for the Web, but papers on all aspects of the Web and Databases are solicited.

Topics of interest include (but are not limited to):

- * Business processes for applications on the Web
- * Data-oriented aspects of Web application development environments
- * Data Models for Web Information Systems
- * Query languages and systems for XML and Web data
- * Semi-structured data management
- * Web Information Extraction
- * Information retrieval in semi-structured data and the Web
- * Data integration over the Web
- * Warehousing of Web data
- * Data synchronization from hand-held devices to the Web
- * Data-intensive applications on the Web
- * Methodologies and tools for Web data publishing
- * Transactions on the Web
- * Web services and distributed computing over the Web
- * Security and integrity issues
- * Web-based distributed data management
- * Semantic Web and reasoning on Web data
- * Web Community Data Management Systems
- * Database Support for Social Web 2.0 applications
- * Multimedia Content Production, Storage and Search

IMPORTANT DATES

Abstract and paper submission deadline: Sun, April 6, 2008 (midnight EST)

Notification of acceptance: Mon, May 12, 2008

Workshop date: Fri, June 13, 2008

SUBMISSION INSTRUCTIONS

Authors are invited to submit original, unpublished research papers that are not being considered for publication in any other forum. Papers should be submitted electronically as PDF files and be formatted using the camera-ready templates available at: <http://www.acm.org/sigs/pubs/proceed/template.html>

The submission site is available at <https://cmt.research.microsoft.com/WEBDB2008/>

Papers submitted cannot exceed six pages in length, including reference and appendix.

Besides regular paper submissions, WebDB 2008 also welcomes the submission of posters and software demonstration proposals, to foster interaction on hot topics and ongoing work. Posters and demo proposals should be 2 pages long when formatted using the ACM style. Demonstration proposals should outline the context and highlights of the software to be presented, and briefly describe the demo scenario. Posters submissions should focus on innovative work related to WebDB's topics of interest; we encourage the joint submission of posters describing new concepts and fundamental results and of demo proposals of software developed based on those new concepts. Electronic versions of the papers will be included in the ACM Digital Library. All of the submissions will be handled electronically. Each paper will be reviewed by at least three members of the program committee.

WORKSHOP CHAIRS:

* Piero Fraternali, Politecnico di Milano, Italy

* Christoph Koch, Cornell University, USA

PROGRAM COMMITTEE

Amer-Yahia Sihem, Yahoo Research, USA
Benedikt Michael, Oxford University, UK
Casati Fabio, Università di Trento, Italy
Cabibbo Luca, Università di Roma 3, Italy
Chan Chee Yong, National University of Singapore, Singapore
Deutsch Alin, UC San Diego, USA
Doan AnHai, University of Wisconsin, USA
Dolog Peter, Aalborg University, Denmark
Dong Luna, AT&T Labs-Research, USA
Dustdar Schahram, Vienna University of Technology, Austria
Gertz Michael, UC Davis, USA
Goh Angela, Nanyang Technological University, Singapore
Houben Geert Jan, Vrije Universiteit Brussel, Belgium
Michel Sebastian, EPFL, Switzerland
Miklau Gerome, University of Massachusetts, USA
Nejd Wolfgang, University of Hannover, Germany
Petropoulos Michalis, University of Buffalo, USA
Scherzinger Stefanie, IBM, Germany
Schewe Klaus-Dieter, Massey University, New Zealand
Shan Ming-Chien, SAP, USA
Shanmugasundaram Jayavel, Yahoo Research, USA
Teniente Ernest, Universidad Politecnica de Catalunya, Spain

XIME - P 2008
5th International Workshop on
XQuery Implementation, Experience and Perspectives

June 13, 2008

<http://www.ximep2008.org/>

Workshop Focus and Theme

XIME-P 2008 invites original research contributions as well as reports on industrial efforts on the implementation, utilization, and overall prospects of XQuery. Like the earlier editions of the XIME-P workshop series, XIME-P 2008 will be held just after ACM SIGMOD/PODS conference (<http://www.sigmod08.org/>), in Vancouver, Canada.

One of the fascinating aspects of XQuery is that work on the language specification and its implementation is happening on the verge of database systems, information retrieval, document processing, and programming languages. For example, XQuery full-text extensions aim at striking a balance between the worlds of structured and unstructured data. XQuery has attracted users from a wide variety of application domains, and its use is not limited to traditional DB server architectures. Computer science research and industry have found quite a number of promising -- and sometimes completely disjoint -- avenues to approach the challenges resulting from these different XQuery usage scenarios. This "heterogeneity" in contributions and attendees has been a source of lively discussions, panels, and an interesting technical program for previous XIME-P editions.

The XIME-P 2008 program will feature talks and panels on research as well as demonstrations and industrial efforts on the implementation and utilization of XQuery.

XIME-P 2008 Topics of Interest

Topics of interest include the following (though interesting and/or innovative papers on all aspects of XQuery are welcome):

- XQuery variants and extensions
 - Coherent XQuery subsets
 - Embedded XQuery
 - Distributed XQuery
 - Scripting /programming with XQuery
 - XQuery full-text extensions
- XQuery applications and architectures
 - The role of XQuery in Web 2.0
 - XQuery and computing in the sciences
 - XQuery for information retrieval
- Optimization of XQuery for demanding applications
 - Compilation vs. interpretation
 - Cost-based optimization strategies
 - Schema-awareness in storage and processing
- XQuery lessons learned
 - Performance evaluation
 - XQuery debugging
 - Teaching XQuery

Keynote Speech

Jim Melton (Oracle, XQuery W3C Working Group Chair) will deliver the XIME-P 2008 workshop's keynote address.

Paper Submission

XIME-P 2008 calls for original contributions relevant to the open list of topics sketched above. We explicitly welcome "war stories", and reports on innovative, off-beat, and "early stage" approaches to the implementation and application of XQuery as long as the submission meets the high quality standards of the XIME-P workshop series.

Papers should not exceed 6 pages in length (including references and appendices), and be submitted in PDF. They should be formatted according to the ACM guidelines and SIG proceedings templates available at <http://www.acm.org/sigs/pubs/proceed/template.html>. More details on the submission process will be given on the XIME-P 2008 web site (<http://www.ximep2008.org>).

Workshop Proceedings Publication

The primary publication medium for XIME-P has been and will be SIGMOD DiSC. This mode of publication ensures wide dissemination and high visibility (e.g., in the ACM Digital Library and Michael Ley's DBLP index). Online proceedings will be additionally hosted at the workshop web site, <http://www.ximep2008.org>. For inclusion in SIGMOD DiSC, we will ask the authors to transfer their copyrights accordingly.

Important Dates

- Paper submission: Fri, March 28, 2008
- Notification of acceptance: Fri, May 2, 2008
- Camera-ready papers due: Fri, May 16, 2008
- Workshop: Fri, June 13, 2008

XIME-P 2008 Workshop Co-Chairs

Carl-Christian Kanne
University of Mannheim
Mannheim, Germany
kanne@informatik.uni-mannheim.de

Fatma Özcan
IBM Almaden Research Center
San Jose, CA, USA
fozcan@almaden.ibm.com

Program Committee

- Andrey Balmin, IBM Almaden Research Center
- Denilson Barbosa, University of Calgary
- Giorgio Ghelli, University of Pisa,
- Torsten Grust, Technical University Munich
- Mary Holstege, MarkLogic
- Yannis Papakonstantinou, UCSD
- Michael Rys, Microsoft
- Jayavel Shanmugasundaram, Yahoo Research
- John Snelson, Oracle
- Jens Teubner, IBM Watson Research Center
- Till Westmann, BEA

Please address questions or comments to the workshop chairs directly.

Estimating the Selectivity of *tf-idf* based Cosine Similarity Predicates

Sandeep Tata Jignesh M. Patel
Department of Electrical Engineering and Computer Science
University of Michigan
2260 Hayward Street, Ann Arbor, Michigan 48109
{tatas, jignesh}@eecs.umich.edu

Abstract

An increasing number of database applications today require sophisticated approximate string matching capabilities. Examples of such application areas include data integration and data cleaning. Cosine similarity has proven to be a robust metric for scoring the similarity between two strings, and it is increasingly being used in complex queries. An immediate challenge faced by current database optimizers is to find accurate and efficient methods for estimating the selectivity of cosine similarity predicates. To the best of our knowledge, there are no known methods for this problem. In this paper, we present the first approach for estimating the selectivity of *tf.idf* based cosine similarity predicates. We evaluate our approach on three different real datasets and show that our method often produces estimates that are within 40% of the actual selectivity.

1 Introduction

A growing number of database applications require approximate string matching predicates on text attributes. For example, in data scrubbing [4] and data integration applications [5, 6], these predicates are valuable in dealing with spelling errors, typographical errors, and problems with non-uniform data representation. Address fields for instance can refer to the same location, but be written using different conventions (“1301 Beal Ave., Ann Arbor” vs. “1301 Beal Avenue, Ann Arbor”). Another example is the case of item descriptions which vary

slightly from vendor to vendor. One might want to search on the description field to find similar items.

For many real world application, the authors in [3, 8] show that the cosine similarity metric can robustly handle spelling errors, rearrangement of words, and other differences in strings. They also demonstrate that cosine similarity searches and joins can be implemented completely in SQL without adding any code to the relational engine. While cosine similarity is a good metric for comparing strings, to the best of our knowledge, there are no known methods for estimating the selectivity of these predicates. As a result, optimizers may often produce inefficient plans for queries involving these predicates. With the increasing use of cosine similarity predicates, there is an urgent need to develop methods that can estimate the selectivity of these predicates.

In this paper, we discuss a technique for estimating the selectivity of *tf.idf* based cosine similarity predicates. We make use of a statistical summary of the distribution of different tokens in the database. We also make use of the distribution of the dot product of a typical query with a database row’s *tf.idf* vector. We present two techniques that use the data in different ways and compare their performance on different datasets.

The rest of the paper is organized as follows: Section 2 describes related work and briefly reviews cosine similarity. Section 3 describes the summary structure we employ. Section 4 describes the algorithm used to compute the estimates. The experimental evaluation is presented in Section 5. Finally, we make concluding remarks and point to directions of future work in Section 6.

2 Review and Related Work

Cosine similarity is a vector-based measure of the similarity of two strings. The basic idea behind cosine similarity is to transform each string into a vector in some high dimensional space such that similar strings are close to each other. The cosine of the angle between two vectors is a measure of how “similar” they are, which in turn, is a measure of the similarity of these strings. If the vectors are of unit length, the cosine of the angle between them is simply the dot product of the vectors.

There are many ways of transforming a string in the database into a vector. The *tf.idf* vector is a popular choice for this representation. The *tf.idf* vector is composed of the product of a *term frequency* and the *inverse document frequency* for each token that appears in the string. The process of constructing the *tf.idf* vector is described below.

As a first step towards implementing the cosine similarity predicate, we construct a *tf.idf* vector for each row in the relation. If there are multiple string attributes of interest in each row, then we need to compute a vector for each string attribute. To keep the discussion simple, we will assume there is only one string attribute in the relation that is used in a cosine similarity operation.

The length of the *tf.idf* vector is equal to the total number of tokens. A token can be a q-gram or a word. If we are using q-grams, then the length of each vector is the total number of possible q-grams = $|A|^q$, where $|A|$ is the size of the alphabet. The vector stores the *tf.idf* value corresponding to each token for each string. The *term frequency* is the number of times the token appears in the string and is a measure of the importance of that token in the string. The *inverse document frequency* (inverse of the number of strings in which the token appears) serves to normalize the effect of tokens (like “the”) that appear commonly in many strings. The product of *tf* and *idf* is a measure of the importance of the token in the string and the database as a whole. Note that in most real datasets, the strings are very short when compared to the total number of possible tokens, and therefore these vectors tend to be very sparse.

When a query comes in, the normalized *tf.idf* vector corresponding to the query is constructed. The *idf* of each term in the query is just 1. We compute the dot product of this vector with the vector for each row in the

database: this is the cosine similarity. If the query and the string share more terms, the dot product is higher. In addition, if they share more “uncommon” terms, that contributes to the score more. The predicate is typically of the form *cosine_similarity*(*R.s*, “*Dr. Jekyll*”) > 0.5, and is evaluated by selecting all those strings where the dot product exceeds the given threshold [3, 8].

To the best of our knowledge, there is no literature on techniques to estimate the selectivity of a cosine similarity predicate. The work closest to ours is [7] where the authors describe a selectivity estimation technique for a fuzzy string predicate. However, this fuzzy predicate is different from any of the well known predicates and has not been shown to perform like cosine similarity in real world tasks [3, 8].

In this paper, we focus on *tf.idf* based cosine similarity. Although there are other vector representations where cosine similarity can be used, *tf.idf* is a popular choice in many applications because of its simplicity and robustness. The techniques in this paper take advantage of some of the properties of *tf.idf*, and therefore will likely require adaptation to work with other vector representations.

Interestingly, the authors of [1] show that many metric distance measures follow a power law distribution for average number of neighbors with respect to distance. That is, number of neighbors within distance s is proportional to s^d where d is some positive constant. As has been argued in [7], this property does not hold for similarity functions like the edit distance, and in our case, the cosine similarity function because of the large number of pairs of words within the same distance. Furthermore, this approach only estimates the average number of neighbors for a string in a dataset, and does not estimate the number of neighbors for a given query string which could be very different from the average.

3 Summary Structure

The summary structure we describe stores a concise representation of the distribution of the *tf.idf* values for each token. If we think of the *tf.idf* vectors for all tuples in the relation as a matrix, we observe that this matrix is very sparse. Table 1 shows a sketch of such a matrix. Most rows in this table are sparse because a given string

rowID	string	Tok 1	Tok N
1	s_1	w_1^1	...	w_N^1
.
.
R	s_R	w_1^R	...	w_N^R
		μ_1, σ_1, C_1	...	μ_N, σ_N, C_N

Table 1: A Table and the *tf.idf* Vectors

is likely to contain only a small number of tokens. In addition, most columns in this table are also sparse, because very few tokens (like “the”) are likely to appear in a large number of strings. We have observed empirically that the probability density function of the *tf.idf* weights for a column is characterized by a large mass of probability at zero (most tokens appear only in a few strings). The rest of the probability is distributed around a small positive value. The proposed summary structure (as shown in the last row of Table 1) captures this distribution by storing the following three values for each token:

1. Mean of X for $X \neq 0$ (μ_i),
2. Standard Deviation for $X \neq 0$ (σ_i), and
3. Probability that a q-gram is non zero, $1 - \text{Prob}(X=0)$ (C_i)

Assume that all the nonzero *tf.idf* values are stored in a table called *Vectors(token,row,value)*. That is, for each token in the original database, the table *Vectors* stores a record for each row in which this token appears with the *tf.idf* value for that token in that row. This is merely a compact way of storing the (sparse) *tf.idf* vector for each row of the database. The summary structure can be generated from *Vectors* by simply using the following SQL query:

```
SELECT token, avg(value),
stddev(value),count(row)/total_rows
FROM Vectors GROUP BY token;
```

Note that the size of this summary structure is bounded by the number of distinct tokens in the language from which the text is drawn. For example, [3, 8] show that for many real applications cosine similarity works well with a token size of three. Assuming an alphabet of size 50 (characters, numbers, punctuation, etc.) the maximum number of tokens (q-grams) is 125K. Also note that the size of the structure is largely independent of the

database size, and for a large text database the summary structure is a very small proportion of the total size.

4 The Estimation Algorithm

The cosine similarity is the dot product of two *tf.idf* vectors representing the query and the database string. The key to estimating the selectivity of a cosine similarity predicate is to understand the distribution of the dot product. In other words, the problem at hand is to compute the cumulative distribution function of the dot product given a) the query vector, and b) a distribution characterizing the *tf.idf* vectors in the database. Once we compute this cumulative probability distribution function, calculating the probability that the cosine similarity exceeds a certain threshold becomes fairly simple.

We model the *tf.idf* vector in the database as a vector of random variables ($X_1 X_2 X_3 \dots X_n$) – one for each token. The dot product can now be modeled as:

$$Y = \sum_{i=1}^n u_i \times X_i \quad (1)$$

where u is the query vector.

A straightforward approach to understanding the distribution of Y is to model the distribution of each of the X_i 's and analytically compute the PDF of Y . However, this turns out to be extremely difficult for any non-trivial characterization of X_i . Alternately if we were to evaluate the PDF of Y by sampling the PDF's of X_i 's, it turns out that the number of samples required to accurately estimate the selectivity is prohibitively high. We therefore choose an alternate technique where we model the distribution of Y and try to determine the parameters of the distribution.

In order to understand how the dot product is distributed, we generated a large set of sample queries by randomly picking strings in the database and introducing one or two errors in the string. We repeated this experiment for a variety of datasets. We observed that the distribution is as shown in Figure 1. The distribution is characterized by a mass of probability close to zero. The rest of the probability is distributed such that it peaks at a small positive value and a long tail tapering off to 0 at $Y = 1$. After evaluating several well known distributions, we determined that this data was modeled accurately as an inverse normal distribution [9] with a mass of proba-

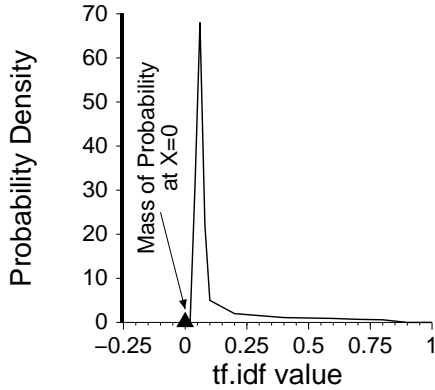


Figure 1: Typical Distribution of the tf.idf Dot Product

bility at 0. We show in the next section, that this simple empirical observation leads to surprisingly good results.

The inverse normal distribution can be completely characterized by its mean and standard deviation. We remind the reader that the probability distribution function of the inverse normal function is:

$$PDF = \sqrt{\frac{B}{2\pi y^3}} \exp\left(-\frac{B}{2y} \left(\frac{y-A}{A}\right)^2\right) \quad (2)$$

where the mean is A , and the variance is $\frac{A^3}{B}$. The cumulative distribution function (CDF) is:

$$CDF = \Phi\left(\sqrt{\frac{B}{y}} \frac{y-A}{A}\right) + \exp\left(\frac{2B}{A}\right) \Phi\left(\sqrt{\frac{B}{y}} \frac{-y-A}{A}\right) \quad (3)$$

where $\Phi(x)$ is the CDF of a standard Gaussian.

The problem now reduces to determining the parameters of the inverse normal distribution (mean and variance). We now present two empirical algorithms for estimating the mean and standard deviation of Y using the data at hand. We then show through experiments in Section 5 that these techniques lead to good estimates.

4.1 Algorithm ES

A simple approach to estimate the mean of Y is to use the weighted average of the means of X_i 's (in Equation 1)

from each column of the *tf.idf* matrix. The mean of each X_i is available in the summary structure. We compute

$$\mu^{ES} = \alpha \times \sum (C_i \times \mu_i \times u_i) \quad (4)$$

where C_i is the probability that token i assumes a nonzero value. μ_i is the mean of the nonzero values of token i as stored in the summary, and u_i is the *tf.idf* weight of the token i in the query vector.

The standard deviation is also computed similarly:

$$\sigma^{ES} = \beta \times \sum (C_i \times \mu_i \times u_i) \quad (5)$$

We call this simple approach ES.

In the above equations, α and β are empirically determined scaling constants for a given relation. They are present to accommodate for the fact that the weighted sum of means does not necessarily yield the actual mean of Y . In order to determine α , we first assume $\alpha = 1$. We determine the average value of the ratio $\frac{\mu^{actual}}{\mu^{ES}}$ for a training set of queries and set α to this value. β is determined similarly. Using samples from a real workload for the training set will ensure that these values are more accurate.

Algorithm ES(query,threshold,summary)

1. Construct the *tf.idf* vector u for the query.
2. Compute $\mu_{ES} = \alpha \sum_{i=1}^N (\mu_i \times C_i \times u_i)$
3. Compute $\sigma_{ES} = \beta \sum_{i=1}^N (\sigma_i \times C_i \times u_i)$
4. Compute over nonzero u_i :
5. $nz_{ES} = 1 - (\prod_{i=1}^N (1 - C_i))^{1/q}$
6. Compute Estimate = $nz_{ES} \times inv_normal_cdf(threshold, \mu_{ES}, \sigma_{ES})$

Figure 2: Estimation using ES

In order to completely characterize Y , we also need to estimate the mass of probability at $Y = 0$. This is the probability that the dot product is zero. We use:

$$PZ^{ES} = (\prod_{i=1}^N (1 - C_i))^{1/q} \quad (6)$$

where N is the number of nonzero *tf.idf* weights in the query vector, and q is the length of the tokens used. In effect, we are computing the product of all the values

corresponding to the nonzero entries in the query vector. The exponentiation with $\frac{1}{q}$ is to correct for the fact that q-grams are usually not independent. For instance tokens like ‘THA’ and ‘HAT’ are more likely to co-occur because they constitute common words like ‘THAT’. This simple approximation leads to some very good estimates.

Once we have μ^{ES} , σ^{ES} , and PZ^{ES} , we calculate the selectivity s of the query as:

$$s = (1 - PZ^{ES}) \times \text{inv_cdf}(\text{threshold}, \mu^{ES}, \sigma^{ES}) \quad (7)$$

where inv_cdf is the CDF for the inverse normal distribution, and threshold is the value obtained from the predicate of the form $\text{cosine_similarity}(\text{R.a}, \text{string}) \geq \text{threshold}$.

4.2 Algorithm EL

Although Algorithm ES gives us fairly good estimates, we found that instead of simply learning constants α and β from a training workload, learning a simple function using linear regression can significantly improve the accuracy of the estimate.

Algorithm EL trains functions to estimate the actual mean and the actual standard deviation for the dot product from μ^{ES} and σ^{ES} computed as in ES using $\alpha = 1$ and $\beta = 1$. In the training phase, we use the data from a set of sample queries that is representative of the workload. We train functions f_μ and f_σ to estimate μ^{actual} and σ^{actual} from μ^{ES} and σ^{ES} . We also train a function to better estimate PZ^{actual} using $PZ_{corrected}^{ES}$. If there are changes to the query workload or the data itself, one can retrain these functions to increase their accuracy. (If such retraining is not feasible, then one can resort to the ES algorithm.) For the training function, we empirically tried and evaluated several families of function, including polynomials of various degrees, exponential functions, and combinations of polynomials and exponentials. We found that the following simple family of functions works best for training the estimators:

$$f(x) = c_1 + c_2x + c_3e^{-x^2} \quad (8)$$

5 Experimental Evaluation

In this section, we present an experimental evaluation of the estimates produced by the ES and EL algorithms on

EstimateEL(query, threshold, summary, f_μ , f_σ , f_{nz})

1. Construct the *tf.idf* vector u for the query.
2. Compute $\mu_{eq} = \sum_{i=1}^N (\mu_i \times C_i \times u_i)$
3. Compute $\sigma_{eq} = \sum_{i=1}^N (\sigma_i \times C_i \times u_i)$
4. Compute over nonzero u_i :
5. $nz_{eq} = 1 - (\prod_{i=1}^N (1 - C_i))^{1/k}$
6. $\mu_{EL} = f_\mu(\mu_{eq})$
7. $\sigma_{EL} = f_\sigma(\sigma_{eq})$
8. $nz_{EL} = f_{nz}(nz_{eq})$
9. Compute Estimate =
 $nz_{EL} \times \text{inv_normal_cdf}(\text{threshold}, \mu_{EL}, \sigma_{EL})$

Figure 3: Estimation using EL

three datasets. The results are largely representative of many other datasets that we tried. The three dataset that we use are SCH, AUT, and HEAD as described below:

1. SCH consists of 99,632 records with high school names and addresses in the USA. The total size of the dataset is 13MB. The school name field was used for cosine similarity.
2. AUT [2] is a set of 371,022 author names from DBLP totaling 8MB.
3. HEAD [10] contains 119,015 article headlines from the Wall Street Journal totaling 7.5MB.

For each dataset, we randomly chose a set of 50 strings from the database itself, and posed 5 queries with each string by varying the cosine similarity threshold from 0.2 to 0.6 in increments of 0.1. Another (different) set was similarly generated to first train ES and EL. Queries were roughly classified as having Low, Medium, or High selectivity based on whether they selected $> 10\%$, $1\% - 10\%$ or $< 1\%$ of the rows respectively.

The size of the summary structure was less than 3% and took less than 3 minutes to construct in each case. Both ES and EL are efficient and take less than 1 millisecond per query to compute an estimate.

We report the average percentage error in Figures 4, 5, and 6. That is, we report $\frac{|\text{estimate} - \text{actual}|}{\text{actual}} \times 100$. The

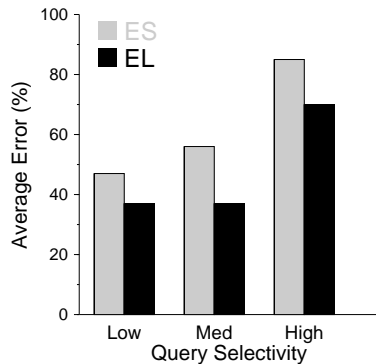


Figure 4: SCH Dataset

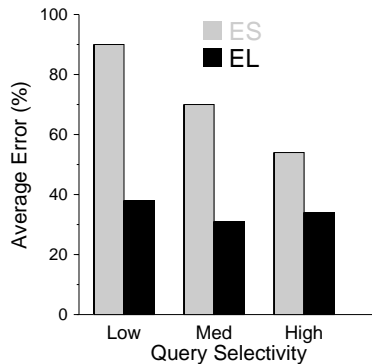


Figure 5: AUT Dataset

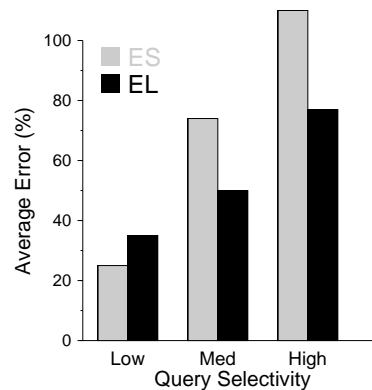


Figure 6: HEAD Dataset

figures show that in each case EL is more accurate than ES by 10 to 50 percentage points. For instance, in Figure 5, in the case of low selectivity queries, ES incurs a 90% error while EL has less than 40% error. Although ES is fairly accurate in many cases, it occasionally has a very large error (eg. high selectivity queries in SCH and HEAD). In all low and medium selectivity cases, the estimates provided by EL have less than 40% error. The error is usually higher in the case of highly selective queries as can be expected. The benefits of using the more complex learning model in EL are evident as they pay off in terms of more accurate estimates.

6 Conclusions and Future Work

In this paper, we have presented the problem of estimating the selectivity of cosine similarity predicates. To our knowledge, this is the first paper to address this problem. We discussed why estimating the selectivity of cosine similarity predicates is a very difficult problem, and proposed a solution based on careful empirical observations about the distribution of the dot product of typical queries. We showed that the approach is space efficient (summaries are small in size) and time efficient (estimation time is also small). We also showed that this technique has reasonably good accuracy in practice.

Directions for future work include exploring analytical modeling for the tf.idf dot product, and alternative approaches that might lead to more accurate estimates.

References

- [1] Caetano Traina and Agma J. M. Traina and Christos Faloutsos. Distance Exponent: A New Concept for Selectivity Estimation in Metric Trees. In *ICDE*, pages 195–195, 2000.
- [2] Digital Bibliography and Library Project (DBLP), <http://dblp.uni-trier.de/>.
- [3] L. Gravano, P. G. Ipeirotis, H. V. Jagadish, N. Koudas, S. Muthukrishnan, L. Pietarinen, and D. Srivastava. Using q-grams in a DBMS for approximate string processing. *IEEE Data Engineering Bulletin*, 24(4):28–34, 2001.
- [4] L. Gravano, P. G. Ipeirotis, N. Koudas, and D. Srivastava. Text Joins for Data Cleansing and Integration in an RDBMS. In *ICDE*, pages 729–731, 2003.
- [5] L. Gravano, P. G. Ipeirotis, N. Koudas, and D. Srivastava. Text Joins in an RDBMS for Web Data Integration. In *WWW*, pages 90–101, 2003.
- [6] Y. Huang and G. Madey. Web Data Integration Using Approximate String Join. In *WWW*, pages 364–365.
- [7] L. Jin and C. Li. Selectivity Estimation for Fuzzy String Predicates in Large Data Sets. In *VLDB*, pages 397–408, 2005.
- [8] N. Koudas, A. Marathe, and D. Srivastava. Flexible String Matching Against Large Databases in Practice. In *VLDB*, pages 1078–1086, 2004.
- [9] A. Papoulis. *Probability, Random Variables and Stochastic Processes*. McGraw-Hill.
- [10] The LDC Corpus Catalog, <http://wave.ldc.upenn.edu/Catalog/>.