26th IEEE International Conference on Data Engineering (ICDE 2010)

March 1–6, 2010 Long Beach, USA

http://www.icde2010.org

Long Beach, USA

TABLE OF CONTENTS

| Message from the ICDF 2010 Chairs | 3 |
|-----------------------------------|-----|
| Vonue Man and Eleornian | 5 |
| | 0 |
| Program at a Glance | 8 |
| Monday | 8 |
| Tuesday | 9 |
| Wednesday | 10 |
| Thursday | 11 |
| Friday | 12 |
| Saturday | 13 |
| Program Details | 14 |
| Monday | 14 |
| Tuesday | 18 |
| Wednesday | 26 |
| Thursday | 32 |
| Friday | 37 |
| Saturday | 43 |
| ICDE 2010 Keynotes | 46 |
| ICDE 2010 Banquet | 48 |
| ICDE 2010 Seminars | 49 |
| ICDE 2010 Paper Abstracts | 51 |
| Research Sessions | 51 |
| Industrial Sessions | 88 |
| Demo Sessions | 91 |
| Conference Co-Located Workshops | 100 |
| Organizing Committees | 112 |

Message from the ICDE 2010 Program Chairs and General Chairs

The 2010 edition of the IEEE International Conference on Data Engineering is being held in the city of Long Beach, California, USA, from March 1st to 6th. This could be viewed as a homecoming for ICDE since the conference was started in Los Angeles in 1984, and Long Beach is a coastal town in the greater Los Angeles area. ICDE was continuously held in Los Angeles until 1990 before starting its world travel to exotic places. During this time, ICDE has established itself as a premier forum in the area of data management, providing a unique opportunity for database researchers, users, practitioners, and developers to exchange new ideas, approaches, and methodologies.

The ICDE 2010 program is composed of the traditional elements: invited keynote talks, research and industrial paper sessions, demonstrations, seminars, panels, and accompanying workshops. We are fortunate to have three distinguished keynote speeches: a joint keynote by Richard Winter of Winter Corporation and Pekka Kostamaa of Teradata, and two keynotes by Jeffrey Naughton of University of Wisconsin-Madison, and Donald Kossmann of ETH-Zurich; they all readily agreed to share their perspectives on the state of database industry and research. In addition, Gio Wiederhold of Stanford University accepted our invitation to be the banquet speaker; Gio was the founder of the ICDE conference series, and we invited him to help us celebrate the return of the conference to the Los Angeles area after a gap of around two decades.

ICDE 2010 received 523 research manuscripts, 18 industrial contributions, and 76 demo proposals, for a total of over 600 submissions. This is a gratifyingly large number considering the global recession of 2009. Moreover, it is a testament to our community's high regard for ICDE, and we thank all authors for submitting their innovative work to the conference.

To facilitate selection of research papers, we organized the program committee into 15 topic-based tracks. Each track was headed by a vice-chair who formed a committee to evaluate the papers assigned to that track. This resulted in a research program committee consisting of almost 230 members in total, with an average review load of around 7 papers. The evaluation process consisted of three distinct phases: initial reviews of the papers by PC members, author responses to these reviews, and finally, PC's discussion and fine-tuning of the reviews. This was followed by a meeting of the vice-chairs and the PC co-chairs during VLDB 2009 in Lyon, France, where the final selections were made.

The research program features 69 long papers and 41 short papers, with long papers presented in 30 minute slots, while short papers have 15 minute slots. The industrial program includes 8 submitted papers and one invited paper. In addition, 29 demonstration proposals were selected for presentation to the conference audience. All research and industrial track papers will be displayed in two plenary poster sessions, a repeat of the successful experiment conducted by ICDE 2009 in Shanghai, China.

Thanks to the diligence of our PC members, more than 95% of the reviews were received on time, and the remaining few were completed shortly afterwards. In a first for ICDE, we set up outstanding vice-chair and outstanding reviewer awards. These awards were intended to recognize both the timeliness and quality of the reviews. For assessing quality, we incorporated mechanisms for authors to rate the thoroughness and fairness of their reviews, and designed simple but informative metrics to characterize the performance of reviewers and vice-chairs. Based on these assessments, Susan Davidson of the University of Pennsylvania received the outstanding vice-chair award, while Sarah Cohen-Boulakia of LRI

France picked up the outstanding reviewer award. They were both granted free conference registrations as a token of our appreciation. We also recognized twenty PC members, whose names appear in the Outstanding Reviewer section of this proceedings, for the consistently high quality of their reviews as indicated by author ratings. We are deeply grateful to all the vice-chairs and PC members for their dedicated and diligent service in coming up with an exciting program that maintains the high standards associated with ICDE.

To our knowledge, ICDE 2010 is the first database conference to provide the authors with an opportunity to rate their reviews and reviewers. They could choose from the following scores: 1 for superficial or biased (Poor), 2 for competent and fair (Satisfactory), and 3 for thoughtful and constructive (Good). The author ratings were generally positive and we are extremely pleased to report that the average score was 2.2, a resounding endorsement of the expertise and professionalism of the program committee.

A committee consisting of the three PC chairs and two vice-chairs (Felix Naumann and Michael Benedikt) had the challenging task of selecting the award-winning research papers from the rich corpus of quality submissions accepted for the conference. "TASM: Top-k Approximate Subtree Matching" by Augsten, Barbosa, Bohlen and Palpanas received the Best Paper award for providing an elegant solution to the classical problem of identifying subtrees in a data tree with the smallest edit distances from a given query tree. "Usher: Improving Data Quality With Dynamic Forms" spearheaded by students from UC Berkeley and MIT received the Best Student Paper award for developing principled machine-learning based techniques to ensure data quality right at its very root: when humans enter data via forms. We congratulate the authors of these papers on their exemplary efforts.

In addition to the paper sessions and demonstrations, the conference program includes five seminars and two panels on topics of current interest to the data engineering community. Accompanying the main conference are seven workshops, four preceding the conference on Monday, March 1, and three following the conference on Friday, March 5, and Saturday, March 6. To provide a forum for student participants to share their thesis research, there is also a PhD Workshop, which partially overlaps the program on Friday, March 5.

The success of ICDE 2010 is a result of collegial teamwork from many individuals, who worked tirelessly to make the conference a top research forum. We acknowledge and thank the sterling endeavors of Mike Carey, Fabio Casati, and Edward Chang (Industrial Chairs), Ioana Manolescu and Sharad Mehrotra (Demo Chairs), Luis Gravano and Sunita Sarawagi (Seminar Chairs), Anastasia Ailamaki and Carlo Zaniolo (Panel Chairs), Christian Jensen and Renee Miller (Workshop Chairs), and Nikos Mamoulis, Yannis Papakonstantinou, and Timos Sellis (PhD Workshop Chairs), and the organizers of the seven accompanying workshops. We also express our deep appreciation of the outstanding work put in over many months by Chen Li (Local Arrangements Chair), Christof Boernhoevd (Finance Chair), Roger Zimmermann (Publicity Chair), and Feifei Li and Mirella Moro (Proceedings Chairs). Without their tireless efforts, this conference would not have been a success. We are also thankful to the many student volunteers from UC Irvine, UCLA, UC Riverside and USC.

In addition, there are many other individuals whose contributions we warmly acknowledge. We benefited greatly from the sage advice provided by Paul Larson, our ICDE Steering Committee Liaison, Calton Pu, ICDE Steering Committee Chair and General Chair of ICDE 2006, and Malu Castellanos, General Chair of ICDE 2008. We also acknowledge the professional help and support provided by Carmen Saliba of the IEEE Computer Society in navigating the many processes involved in organizing an IEEE conference. The 2009 Program Chairs, Yannis Ioannidis, Dik Lee, and Raymond Ng, continually provided the PC Chairs with valuable tips and

we are grateful to them. We express our gratitude to the Microsoft CMT team for their ready assistance and quick replies to our multitude of requests. We thank the Local Arrangements Chairs of VLDB 2009, Mohand-Said Hacid and Jean-Marc Petit, who generously helped with the logistics of the PC meeting in Lyon, France.

We warmly acknowledge the financial support of our corporate sponsors: Microsoft (at the platinum level); HP, IBM and Greenplum (at the silver level); and AT&T, ESRI, Google, Oracle, SAP, and Yahoo! Labs (at the bronze level). In addition, we thank Maria Zemankova and Frank Olken of NSF for a generous grant that enabled us to provide travel support for 27 Ph.D. students from various US universities to attend the conference.

Finally, we thank all the authors, presenters, and participants of the conference. We hope all of you enjoy the conference and the surrounding Long Beach area! This edition of ICDE shows that data engineering continues to be a thriving and vibrant discipline!

| Shahram Ghandeharizadeh | Umeshwar Dayal |
|-------------------------|------------------|
| Jayant Haritsa | Vassilis Tsotras |
| Gerhard Weikum | |

ICDE 2010 PC Chairs

ICDE 2010 General Chairs

VENUE MAP AND FLOORPLAN



Hyatt Regency Long Beach

DIRECTIONS

From LAX: Take 405 South to 710 South. Follow signs to Convention Center / Shoreline Drive. Turn left on Pine. From Orange County / John Wayne Airport: Take 405 North to 710 South. Follow signs to Convention Center / Shoreline Drive. Turn left on Pine. From Long Beach Airport: Exit right on Lakewood Blvd. Turn right on Spring St. to 405 North, then to 710 South. Follow signs to Convention Center. Turn left on Pine Ave.

> Hyatt Regency Long Beach 200 South Pine Avenue, Long Beach, California, USA 90802 Tel: +1 562 491 1234 Fax: +1 562 983 1491



ICDE 2010

| | Program at a Glance — Monday, March 1st | | | | | | | | | | | |
|---------------|---|---|--|---|--|--|--|--|--|--|--|--|
| When | Regency A | A Regency B Regency C Regency D Regency EF | | Beacon A | Beacon B | | | | | | | |
| 08:00-08:30 | Continental Breakfast (Regency Foyer) | | | | | | | | | | | |
| 08:30 - 10:00 | | DBRank Opening Remarks and Keynote: Alon Halevy (Google) | MOUND Opening Remarks and Keynote: Dan Olteanu (Oxford University) | WISS Opening Remarks and Keynote I: Mike Carey (UC Irvine) (08:45 - 10:00) | SMDB Opening Remarks and Paper Session I (Full papers, Indexing and Workload Management) | | | | | | | |
| 10:00 - 10:30 | | | Coffee I | Break (Regency Foyer) | | | | | | | | |
| 10:30 - 12:00 | .0:30 - 12:00 | | MOUND Paper Session | WISS Research Papers | SMDB Keynote: Oliver Ratzesberger (eBay) | | | | | | | |
| 12:00 - 13:30 | | | Lur | ich (on your own) | | | | | | | | |
| 13:30 - 15:00 | | DBRank Research Session 2 (1:30-2:30) | | WISS Keynote II: Jeff Hammerbacher (Cloudera) | SMDB Paper Session II (Full papers, On Cloud and Column Stores) | | | | | | | |
| 15:00 - 15:30 | | | Coffee I | Break (Regency Foyer) | | | | | | | | |
| 15:30 - 16:30 | | | | WISS Industry Solutions | SMDB Paper Session III (Short papers, Query Optimization and Workload Management) | | | | | | | |
| 16:30 - 18:00 | | | | (15:30 - 17:30) | SMDB Panel (Databases, MapReduce and the Cloud-Oh My! What's in it for the Administrator? | | | | | | | |

Long Beach, USA

 ∞

ICDE 2010

| Program at a Glance — Tuesday, March 2nd | | | | | | | | | |
|--|-----------|--|--|--|-----------------------------------|---|---|--|--|
| When | Regency A | Regency B | Regency C | Regency D | Regency EF | Beacon A | Beacon B | | |
| 08:00-08:30 | | | Contine | ntal Breakfast (Reg | ency Foyer) | | | | |
| 08:30 - 10:00 | | | | | | Opening Ses Richard Wi Kostama Large Warehousin Obse | ssion; Keynote 1 nter and Pekka a (presenter) Scale Data ng: Trends and rvations | | |
| 10:00 - 10:30 | | | Cof | fee Break (Regency | v Foyer) | • | | | |
| 10:30 - 12:00 | Demo 1A | Research 1 KNN Queries | Research 2 Distributed Data | Research 3 Stream Mining | Industry 1 Data Warehousing | | | | |
| 12:00 - 13:30 | | | | | | Lunch (by | v conference) | | |
| 13:30 - 15:00 | Demo 2A | Research 4 Location Based Services | Research 5 Probabilistic Databases | Research 6 Spatial Indexing | Seminar 1 | | | | |
| 15:00 - 15:30 | | | Cof | fee Break (Regency | v Foyer) | | | | |
| 15:30 - 17:00 | | Research 7 Privacy Techniques | Research 8 Skyline Queries | Research 9 Information Integration | Seminar 1 (cont.) | Research 10 Query Interfaces | | | |
| 17:00 - 18:00 | | | | | | Poste | Poster Set-Up | | |
| 18:00 - 20:00 | | | | | | "All Posters" Ses | sion and Reception | | |

9

| _ |
|---------|
| \cap |
| |
| Ĕ |
| ь. • |
| 2 |
| \leq |
| 5 |

| | Program at a Glance — Wednesday, March 3rd | | | | | | | | | |
|---------------|--|--|---|--|--|--|---|--|--|--|
| When | Regency A | Regency B | Regency C | Regency D | Regency EF | Beacon A | Beacon B | | | |
| 08:00-08:30 | | | Contine | ental Breakfast (Re | gency Foyer) | | | | | |
| 08:30 - 10:00 | | | | | | Ke Jeffrey DBMS: Lesson Years, Speculat | ynote 2 F. Naughton s from the First 50 ions for the Next 50 | | | |
| 10:00 - 10:30 | | | Coi | ffee Break (Regenc | y Foyer) | | | | | |
| 10:30 - 12:00 | Demo 1B | Research 11 Top-K Queries | Research 12 Workflow and Workload Management | Research 13 Indexing and Hashing | Industry 2 Data, Data, and More Data | | | | | |
| 12:00 - 14:00 | | | | | | Busin | less Lunch | | | |
| 14:00 - 15:30 | Demo 2B | Research 14 Scientific Data Mining | Research 15 Database Performance and Reliability | Research 16 Spatial Databases | Seminar 2 | | | | | |
| 15:30 - 16:00 | Coffee Break (Regency Foyer) | | | | | | | | | |
| 16:00 - 17:30 | | Research 17 Sensor Networks | Research 18 Query Optimization | Research 19 Graph Mining | Seminar 2 (cont.) | Research 20 Parallel Processing | | | | |

10

ICDE 2010

| Program at a Glance — Thursday, March 4th | | | | | | | | | |
|---|-----------|--|---|---|--------------------|--------------------------|--------------------------------------|--|--|
| When | Regency A | Regency B | Regency C | Regency D | Regency EF | Beacon A | Beacon B | | |
| 08:00-08:30 | | | Contine | ental Breakfast (Reg | gency Foyer) | | | | |
| 08:30 - 10:00 | | | | | | Key Donald How new | note 3: Kossmann is the cloud? | | |
| 10:00 - 10:30 | | | Со | ffee Break (Regency | y Foyer) | | | | |
| 10:30 - 12:00 | | Research 21 Keyword Search | Research 22 Query Processing | Research 23 Web and Collaborative Applications | Panel 1 | Seminar 3 | | | |
| 12:00 - 13:30 | | • | • | Lunch (on your o | wn) | | | | |
| 13:30 - 14:00 | | | | | | Poster | Preparation | | |
| 14:00 - 15:30 | | | | | | "All Pos | ters" Session | | |
| 15:30 - 16:00 | | | Со | ffee Break (Regency | y Foyer) | | | | |
| 16:00 - 17:30 | | Research 24 Scientific Databases | Research 25 Tree Queries and Semi- Structured Databases | Research 26 Query Ranking and Database Testing | Seminar 4 | Panel 2 | | | |
| 18:00 - 19:00 | | Pre-b | oanquet Jazz Prese | ntation by Double | Blind Revue (Reger | ncy Foyer) | | | |
| 19:00 - 21:00 | | | | | | Ва | nquet | | |

11

| ī | |
|----|--|
| DE | |
| 20 | |
| 0 | |

| ICD | | Program at a Glance — Friday, March 5th | | | | | | | | |
|--|---------------|---|--|-------------------------------------|-------------------------------------|--|-------------------|----------|--|--|
| E 20 | When | Regency A | Regency B | Regency C | Regency D | Regency EF | Beacon A | Beacon B | | |
| 10 | 08:00-08:30 | Continental Breakfast (Regency Foyer) | | | | | | | | |
| | 08:30 - 10:00 | | Research 27 Social Networks and Similarity Queries | Research 28 Stream Processing | Industry 3 Query Optimization | Seminar 5 | PhD Workshop 1 | | | |
| 10:00 - 10:30 Coffee Break (Regency Foyer) | | | | | | | | | | |
| - | 10:30 - 12:00 | | Research 29 Publishing Privacy | Research 30 Data Clouds | | Seminar 5 (cont.) | PhD Workshop 2 | | | |
| | 12:00 - 13:30 | Lunch (on your own) | | | | | | | | |
| 12 | 13:30 - 15:00 | | DESWeb Keynote: Howard Ho (IBM Almaden) & Panel | | | NTII Opening Remarks and Keynote 1: Dan Wolfson (IBM) | PhD Workshop 3 | | | |
| ĺ | 15:00 - 15:30 | Coffee Break (Regency Foyer) | | | | | | | | |
| E | 15:30 - 17:00 | | DESWeb Session 1 Data Search | | | NTII Session 1 New Architectures and Models | | | | |
|)ng Beach, U | 17:00 - 18:00 | | DESWeb Session 2 Techniques for Mappings | | | | | | | |

Long Beach, USA

| | Program at a Glance — Saturday, March 6th | | | | | | | | | | |
|---------------|---|---|----------------------------------|---------------------|---|----------|----------|--|--|--|--|
| When | Regency A | Regency B | Regency C | Regency D | Regency EF | Beacon A | Beacon B | | | | |
| 08:00-08:30 | | • | Contine | ntal Breakfast (Reg | gency Foyer) | | | | | | |
| 08:30 - 10:00 | | DESWeb Session 3 Data Management (9:10-10:00) | M3SN Keynote and Session 1 | | NTII Keynote 2 Craig Knoblock (USC) | | | | | | |
| 10:00 - 10:30 | | • | Cof | fee Break (Regency | v Foyer) | | | | | | |
| 10:30 - 12:00 | | DESWeb Session 4 Entity Analysis | M3SN Session 2 | | NTII Session 2 Applications | | | | | | |
| 12:00 - 13:30 | | Lunch (on your own) | | | | | | | | | |
| 13:30 - 15:00 | | | | | NTII Session 3 and Wrap-Up New Primitives and Techniques | | | | | | |

ICDE 2010

PROGRAM DETAILS

08:00-08:30 Continental Breakfast (Regency Foyer)

MONDAY 08:30-10:00

DBRank Opening Remarks and Keynote

Regency B, 08:30 - 10:00, Monday, Workshop

Keynote Speaker: Alon Halevy (Google)

Keynote Title: Table Search

MOUND Opening Remarks and Keynote

Regency C, 08:30 – 10:00, Monday, Workshop

Keynote Speaker: Dan Olteanu (Oxford University)

Keynote Title: A Toolbox of Query Evaluation Techniques for Probabilistic Databases

WISS Opening Remarks and Keynote I

Regency D, 08:30 – 10:00, Monday, Workshop

Keynote Speaker: Mike Carey (UC Irvine)

Keynote Title: Data Services: Past, Present, and Future

SMDB Opening Remarks and Paper Session I: Indexing and Workload Management

Regency EF, 08:30 – 10:00, Monday, Workshop

Adaptive Indexing for Relational Keys

Goetz Graefe, Harumi Kuno; HP, USA

On the Use of Query-Driven XML Auto-Indexing

Karsten Schmidt, Theo Härder; University of Kaiserslautern, Germany

Autonomic Workload Execution Control Using Throttling

*Wendy Powley*¹, *Patrick Martin*¹, *Mingyi Zhang*¹, *Paul Bird*², *Keith McDonald*²; ¹*Queen's University, Canada;* ²*IBM, Canada*

10:00 – 10:30 Coffee Break (Regency Foyer)

MONDAY 10:30-12:00

DBRank Session 1

Regency B, 10:30 - 12:00, Monday, Workshop

Subspace Similarity Search Using the Ideas of Ranking and Top-k Retrieval

Thomas Bernecker, Tobias Emrich, Franz Graf, Hans-Peter Kriegel, Peer Kröger, Matthias Renz, Erich Schubert, Arthur Zimek; LMU München, Germany

Efficient k-Nearest Neighbor Queries with the Signature Quadratic Form Distance

Christian Beecks, Merih Seran Uysal, Thomas Seidl; RWTH Aachen University, Germany

Top-k Pipe Join

Davide Martinenghi, Marco Tagliasacchi; Politecnico di Milano, Italy

MOUND Paper Session

Regency C, 10:30 – 12:00, Monday, Workshop

Constrained Frequent Itemset Mining from Uncertain Data Streams

Carson Kai-Sang Leung, Boyu Hao, Fan Jiang; University of Manitoba, Canada

Cleansing Uncertain Databases Leveraging Aggregate Constraints

*Haiquan Chen*¹, *Wei-Shinn Ku*¹, *Haixun Wang*²; ¹*Auburn University, USA*; ²*Microsoft, China*

U-DBSCAN : A Density-Based Clustering Algorithm for Uncertain Objects

Apinya Tepwankul, Songrit Maneewongwattana; King Mongkut's University of Technology Thonburi, Thailand

WISS Research Papers

Regency D, 10:30 - 12:00, Monday, Workshop

Towards Enterprise Software as a Service in the Cloud

*Jan Schaffner*¹, *Dean Jacobs*², *Benjamin Eckart*¹, *Jan Brunnert*¹, *Alexander Zeier*¹; ¹*HPI, Germany*; ²*SAP, Germany*

The HiBench Benchmark Suite: Characterization of the MapReduce-Based Data Analysis

Shengsheng Huang, Jie Huang, Jinquan Dai, Tao Xie, Bo Huang; Intel, China

End-to-End Confidentiality for a Message Warehousing Service Using Identity-Based Encryption

Yuecel Karabulut¹, Harald Weppner¹, Ike Nassi¹, Anusha Nagarajan², Yash Shroff², Nishant Dubey², Tyelisa Shields²; ¹SAP, USA; ²Carnegie Mellon University, USA

SMDB Keynote

Regency EF, 10:30 – 12:00, Monday, Workshop

Keynote Speaker: Oliver Ratzesberger (eBay)

Keynote Title: TBA

12:00-13:30 Lunch (on your own)

MONDAY 13:30-15:00

DBRank Session 2

Regency B, 13:30 – 14:30, Monday, Workshop

On Novelty in Publish/Subscribe Delivery

Dimitris Souravlias, Marina Drosou, Kostas Stefanidis, Evaggelia Pitoura; University of Ioannina, Greece

Ranking for Data Repairs

Mohamed Yakout, Ahmed K. Elmagarmid, Jennifer Neville; Purdue University, USA

WISS Keynote II

Regency D, 13:30 - 15:00, Monday, Workshop

Keynote Speaker: Jeff Hammerbacher (Cloudera)

Keynote Title: Open Questions for Building an Enterprise Data Platform on the Cloud

SMDB Paper Session II: On Cloud and Column Stores

Regency EF, 13:30 – 15:00, Monday, Workshop

Statistics-Driven Workload Modeling for the Cloud

Archana Ganapathi, Yanpei Chen, Armando Fox, Randy Katz, David Patterson; University of California at Berkeley, USA

A Generic Auto-Provisioning Framework for Cloud Databases

*Jennie Rogers*¹, *Olga Papaemmanouil*², *Ugur Cetintemel*¹; ¹*Brown University, USA*; ²*Brandeis University, USA*

Vertical Partitioning of Relational OLTP Databases Using Integer Programming

Rasmus Resen Amossen; IT University of Copenhagen, Denmark

15:00-15:30 Coffee Break (Regency Foyer)

MONDAY 15:30-16:30

WISS Industry Solutions

Regency D, 15:30 - 16:30, Monday, Workshop

SMDB Paper Session III: Query Optimization and Workload Management

Regency EF, 15:30 - 16:30, Monday, Workshop

Caching All Plans with Just One Optimizer Call

Debabrata Dash, Ioannis Alagiannis, Cristina Maier, Anastasia Ailamaki; EPFL, Switzerland

Towards Workload-Aware Self-Management: Predicting Significant Workload Shifts

Marc Holze, Ali Haschimi, Norbert Ritter; University of Hamburg, Germany

Automatic Tuning of the Multiprogramming Level in Sybase SQL Anywhere

*Mohammed Abouzour*¹, *Kenneth Salem*², *Peter Bumbulis*¹; ¹*Sybase iAnywhere, Canada;* ²*University of Waterloo, Canada*

WISS Industry Solutions (cont.)

Regency D, 16:30 - 18:00, Monday, Workshop

SMDB Panel (Databases, MapReduce and the Cloud-Oh My! What's in it for the Administrator?)

Regency EF, 16:30 – 18:00, Monday, Workshop

Panelists:

- · Ashraf Aboulnaga (Univ. of Waterloo)
- Namit Jain (Facebook)
- Guy Lohman (IBM)
- · Oliver Ratzesberger (eBay)
- Benjamin Reed (Yahoo!)
- Jingren Zhou (Microsoft)

08:00-08:30 Continental Breakfast (Regency Foyer)

TUESDAY 08:30 - 10:00

Opening Session; Keynote 1

Beacon AB, 08:30 - 10:00, Tuesday Chair: Shahram Ghandeharizadeh

Large Scale Data Warehousing: Trends and Observations

*Richard Winter*¹, *Pekka Kostamaa*²; ¹*Winter Corporation, USA*; ²*Teradata, USA*

10:00 - 10:30 Coffee Break (Regency Foyer)

TUESDAY 10:30-12:00

Demo Session 1A: Events, Streams, Services, Mashups and Search

Regency A, 10:30 – 12:00, Tuesday

A Demonstration of the MaxStream Federated Stream Processing System

Irina Botan¹, Younggoo Cho², Roozbeh Derakhshan¹, Nihal Dindar¹, Ankush Gupta¹, Laura M. Haas³, Kihong Kim², Chulwon Lee², Girish Mundada⁴, Ming-Chien Shan⁴, Nesime Tatbul¹, Ying Yan⁵, Beomjin Yun², Jin Zhang⁵; ¹ETH Zürich, Switzerland; ²SAP, Korea; ³IBM, USA; ⁴SAP, USA; ⁵SAP, China

E-Cube: Multi-Dimensional Event Sequence Processing Using Concept and Pattern Hierarchies

Mo Liu¹, Elke A. Rundensteiner¹, Kara Greenfield¹, Chetan Gupta², Song Wang², Ismail Ari³, Abhay Mehta²; ¹Worcester Polytechnic Institute, USA; ²HP, USA; ³Ozyegin University, Turkey

TargetSearch: A Ranking Friendly XML Keyword Search Engine

Ziyang Liu, Yichuan Cai, Yi Chen; Arizona State University, USA

Efficient Fuzzy Type-Ahead Search in TASTIER

*Guoliang Li*¹, *Shengyue Ji*², *Chen Li*², *Jiannan Wang*¹, *Jianhua Feng*¹; ¹*Tsinghua University, China;* ²*University of California at Irvine, USA*

MASS: A Multi-Facet Domain-Specific Influential Blogger Mining System

Yichuan Cai, Yi Chen; Arizona State University, USA

Product EntityCube: A Recommendation and Navigation System for **Product Search**

*Jongwuk Lee*¹, *Seung-won Hwang*¹, *Zaiqing Nie*², *Ji-Rong Wen*²; ¹*POSTECH, Korea*; ²*Microsoft, China*

Navigating Through Mashed-Up Applications with COMPASS

Daniel Deutch, Ohad Greenshpan, Tova Milo; Tel-Aviv University, Israel

GenerIE: Information Extraction Using Database Queries

Luis Tari¹, Phan Huy Tu¹, Jörg Hakenberg¹, Yi Chen¹, Tran Cao Son², Graciela Gonzalez¹, Chitta Baral¹; ¹Arizona State University, USA; ²New Mexico State University, USA

Power-Aware Data Analysis in Sensor Networks

Daniel Klan¹, Katja Hose¹, Marcel Karnstedt², Kai-Uwe Sattler¹; ¹Ilmenau University of Technology, Germany; ²NUI Galway, Ireland

A View-Based Monitoring for Privacy-Aware Web Services

Hassina Meziane¹, Salima Benbernou¹, Aouda K. Zerdali¹, Mohand-Said Hacid², Mike Papazoglou³; ¹Université Paris Descartes, France; ²Université de Lyon, France; ³Tilburg University, The Netherlands

Viewing a World of Annotations Through AnnoVIP

Konstantinos Karanasos, Spyros Zoupanos; INRIA, France

MashRank: Towards Uncertainty-Aware and Rank-Aware Mashups

Mohamed A. Soliman, Mina Saleeb, Ihab F. Ilyas; University of Waterloo, Canada

T-Warehouse: Visual OLAP Analysis on Trajectory Data

Luca Leonardi¹, Gerasimos Marketos², Elias Frentzos², Nikos Giatrakos², Salvatore Orlando¹, Nikos Pelekis², Alessandra Raffaetà¹, Alessandro Roncato¹, Claudio Silvestri¹, Yannis Theodoridis²; ¹Università Ca' Foscari Venezia, Italy; ²University of Piraeus, Greece

WikiAnalytics: Ad-Hoc Querying of Highly Heterogeneous Structured Data

Andrey Balmin¹, Emiran Curtmola²; ¹IBM, USA; ²University of California at San Diego, USA

SMARTINT: A System for Answering Queries Over Web Databases Using Attribute Dependencies

Ravi Gummadi, Anupam Khulbe, Aravind Kalavagattu, Sanil Salvi, Subbarao Kambhampati; Arizona State University, USA

Research Session 1: KNN Queries

Regency B, 10:30 – 12:00, Tuesday Chair: Julia Stoyanovich

K Nearest Neighbor Queries and KNN-Joins in Large Relational Databases (Almost) for Free

Bin Yao, Feifei Li, Piyush Kumar; Florida State University, USA

Quantile-Based KNN Over Multi-Valued Objects

Wenjie Zhang, Xuemin Lin, Muhammad Aamir Cheema, Ying Zhang, Wei Wang; University of New South Wales, Australia

Efficient Rank Based KNN Query Processing Over Uncertain Data

Ying Zhang, Xuemin Lin, Gaoping Zhu, Wenjie Zhang, Qianlu Lin; University of New South Wales, Australia

Research Session 2: Distributed Data

Regency C, 10:30 – 12:00, Tuesday Chair: Hank Korth

Reliable Storage and Querying for Collaborative Data Sharing Systems

Nicholas E. Taylor, Zachary G. Ives; University of Pennsylvania, USA

Strongly Consistent Replication for a Bargain

Konstantinos Krikellas¹, Sameh Elnikety², Zografoula Vagena³, Orion Hodson²; ¹University of Edinburgh, UK; ²Microsoft, UK; ³Concentra Consulting Ltd., UK

Detecting Inconsistencies in Distributed Data

Wenfei Fan, Floris Geerts, Shuai Ma, Heiko Müller; University of Edinburgh, UK

Research Session 3: Stream Mining

Regency D, 10:30 – 12:00, Tuesday Chair: Felix Naumann

Optimal Load Shedding with Aggregates and Mining Queries

Barzan Mozafari, Carlo Zaniolo; University of California at Los Angeles, USA

Scheduling for Fast Response Multi-Pattern Matching Over Streaming Events

*Ying Yan*¹, *Jin Zhang*¹, *Ming-Chien Shan*²; ¹*SAP, China*; ²*SAP, USA*

Discovery of Cross-Similarity in Data Streams (Short paper)

Machiko Toyoda, Yasushi Sakurai; NTT, Japan

Mining Distribution Change in Stock Order Streams (Short paper)

*Xiaoyan Liu*¹, *Xindong Wu*², *Huaiqing Wang*³, *Rui Zhang*¹, *James Bailey*¹, *Kotagiri Ramamohanarao*¹; ¹University of Melbourne, Australia; ²Hefei University of Technology, China; ³City University of Hong Kong, China

Industry Session 1: Data Warehousing

Regency EF, 10:30 – 12:00, Tuesday Chair: Paul Larson

Hive — A Petabyte Scale Data Warehouse Using Hadoop

Ashish Thusoo, Joydeep Sen Sarma, Namit Jain, Zheng Shao, Prasad Chakka, Ning Zhang, Suresh Antony, Hao Liu, Raghotham Murthy; Facebook, USA

Tuning Servers, Storage and Database for Energy Efficient Data Warehouses

Meikel Poess¹, Raghunath Othayoth Nambiar²; ¹Oracle, USA; ²HP, USA

A New Algorithm for Small-Large Table Outer Joins in Parallel DBMS

Yu Xu, Pekka Kostamaa; Teradata, USA

12:00 - 13:30 Lunch (by conference, Beacon AB)

TUESDAY 13:30-15:00

Demo Session 2A: Scalability, Design, Optimization and Visualization

Regency A, 13:30 – 15:00, Tuesday

Mini-Me: A Min-Repro System for Database Software

Nicolas Bruno, Rimma V. Nehme; Microsoft, USA

I/O-Efficient Statistical Computing with RIOT

Yi Zhang, Weiping Zhang, Jun Yang; Duke University, USA

Interactive Physical Design Tuning

Nicolas Bruno, Surajit Chaudhuri; Microsoft, USA

Visualizing Cost-Based XQuery Optimization

Andreas M. Weiner, Theo Härder, Renato Oliveira da Silva; University of Kaiserslautern, Germany

XML Reasoning Made Practical Pierre Genevès¹, Nabil Layaïda²; ¹CNRS, France; ²INRIA, France

TransScale: Scalability Transformations for Declarative Applications

Alexander Böhm, Erich Marth, Carl-Christian Kanne; University of Mannheim, Germany

Reverse Engineering Models from Databases to Bootstrap Application Development

Ankit Malpani¹, Philip A. Bernstein², Sergey Melnik³, James F. Terwilliger²; ¹IIT Madras, India; ²Microsoft, USA; ³Google, USA

ICDE 2010

Long Beach, USA

HECATAEUS: Regulating Schema Evolution

*George Papastefanatos*¹, *Panos Vassiliadis*², *Alkis Simitsis*³, *Yannis Vassiliou*¹; ¹*National Technical University of Athens, Greece;* ²*University of Ioannina, Greece;* ³*HP, USA*

ROX: The Robustness of a Run-Time XQuery Optimizer Against Correlated Data

*Riham Abdel Kader*¹, *Peter A. Boncz*², *Stefan Manegold*², *Maurice van Keulen*¹; ¹*University of Twente, The Netherlands*; ²*CWI, The Netherlands*

Symphony: A Platform for Search-Driven Applications

John C. Shafer, Rakesh Agrawal, Hady W. Lauw; Microsoft, USA

ProbClean: A Probabilistic Duplicate Detection System

George Beskales, Mohamed A. Soliman, Ihab F. Ilyas, Shai Ben-David, Yubin Kim; University of Waterloo, Canada

TransDec: A Spatiotemporal Query Processing Framework for Transportation Systems

Ugur Demiryurek, Farnoush Banaei-Kashani, Cyrus Shahabi; University of Southern California, USA

Provenance Browser: Displaying and Querying Scientific Workflow Provenance Graphs

Manish Kumar Anand¹, Shawn Bowers², Bertram Ludäscher¹; ¹University of California at Davis, USA; ²Gonzaga University, USA

Inconsistency Resolution in Online Databases

Yannis Katsis¹, Alin Deutsch¹, Yannis Papakonstantinou¹, Vasilis Vassalos²; ¹University of California at San Diego, USA; ²Athens University of Economics & Business, Greece

Research Session 4: Location Based Services

Regency B, 13:30 – 15:00, Tuesday Chair: Mohamed Mokbel

TrajStore: An Adaptive Storage System for Very Large Trajectory Data Sets

Philippe Cudre-Mauroux, Eugene Wu, Samuel R. Madden; MIT, USA

C3: Concurrency Control on Continuous Queries Over Moving Objects *Jing Dai, Chang-Tien Lu; Virginia Tech, USA*

Policy-Aware Sender Anonymity in Location Based Services

*Alin Deutsch*¹, *Richard Hull*², *Avinash Vyas*³, *Kevin Keliang Zhao*¹; ¹*University of California at San Diego, USA*; ²*IBM, USA*; ³*Bell Labs Research, USA*

Research Session 5: Probabilistic Databases

Regency C, 13:30 – 15:00, Tuesday Chair: Reynold Chang

Approximate Confidence Computation in Probabilistic Databases

Dan Olteanu¹, Jiewen Huang¹, Christoph Koch²; ¹University of Oxford, UK; ²Cornell University, USA

PIP: A Database System for Great and Small Expectations *Oliver Kennedy, Christoph Koch; Cornell University, USA*

Generator-Recognizer Networks: A Unified Approach to Probabilistic Databases (Short paper)

Ruiwen Chen, Yongyi Mao, Iluju Kiringa; University of Ottawa, Canada

Probabilistic Declarative Information Extraction (Short paper)

*Daisy Zhe Wang*¹, *Eirinaios Michelakis*¹, *Michael J. Franklin*¹, *Minos Garofalakis*², *Joseph M. Hellerstein*¹; ¹*University of California at Berkeley, USA*; ²*Technical University of Crete, Greece*

Research Session 6: Spatial Indexing

Regency D, 13:30 – 15:00, Tuesday Chair: Cyrus Shahabi

PARINET: A Tunable Access Method for In-Network Trajectories

*Iulian Sandu Popa*¹, *Karine Zeitouni*¹, *Vincent Oria*², *Dominique Barth*¹, *Sandrine Vial*¹; ¹*PRiSM, France;* ²*New Jersey Institute of Technology, USA*

Multi-Guarded Safe Zone: An Effective Technique to Monitor Moving Circular Range Queries

*Muhammad Aamir Cheema*¹, *Ljiljana Brankovic*², *Xuemin Lin*¹, *Wenjie Zhang*¹, *Wei Wang*¹; ¹*University of New South Wales, Australia;* ²*University of Newcastle, Australia*

Geotagging with Local Lexicons to Build Indexes for Textually-Specified Spatial Data

Michael D. Lieberman, Hanan Samet, Jagan Sankaranarayanan; University of Maryland at College Park, USA

Seminar 1

Regency EF, 13:30 - 15:00, Tuesday

Anonymized Data: Generation, Models, Usage

Graham Cormode, Divesh Srivastava; AT&T Labs Research, USA

15:00-15:30 Coffee Break (Regency Foyer)

TUESDAY 15:30-17:00

Research Session 7: Privacy Techniques

Regency B, 15:30 – 17:00, Tuesday Chair: Wei-Shinn Ku

On Optimal Anonymization for *l*⁺-Diversity

Junqiang Liu, Ke Wang; Simon Fraser University, Canada

Differential Privacy via Wavelet Transforms

*Xiaokui Xiao*¹, *Guozhang Wang*², *Johannes Gehrke*²; ¹*Nanyang Technological University, Singapore;* ²*Cornell University, USA*

Efficient Verification of Shortest Path Search via Authenticated Hints

*Man Lung Yiu*¹, *Yimin Lin*², *Kyriakos Mouratidis*²; ¹*Hong Kong Polytechnic University, China*; ²*Singapore Management University, Singapore*

Research Session 8: Skyline Queries

Regency C, 15:30 – 17:00, Tuesday Chair: Xuemin Lin

Evaluating Skylines in the Presence of Equijoins

*Wen Jin*¹, *Michael D. Morse*¹, *Jignesh M. Patel*², *Martin Ester*³, *Zengjian Hu*³; ¹*University of Michigan, USA*; ²*University of Wisconsin-Madison, USA*; ³*Simon Fraser University, Canada*

Route Skyline Queries: A Multi-Preference Path Planning Approach *Hans-Peter Kriegel, Matthias Renz, Matthias Schubert; LMU München, Germany*

Probabilistic Contextual Skylines

Dimitris Sacharidis, Anastasios Arvanitis, Timos Sellis; Athena RC, Greece

Research Session 9: Information Integration

Regency D, 15:30 - 17:00, Tuesday Chair: Mourad Ouzzani

Schema Covering: A Step Towards Enabling Reuse in Information Integration

*Barna Saha*¹, *Ioana Stanoi*², *Kenneth L. Clarkson*²; ¹*University of Maryland at College Park, USA*; ²*IBM, USA*

Managing Uncertainty of XML Schema Matching

Reynold Cheng, Jian Gong, David W. Cheung; University of Hong Kong, China

Propagating Updates Through XML Views Using Lineage Tracing

Leonidas Fegaras; University of Texas at Arlington, USA

Seminar 1 (cont.)

Regency EF, 15:30 – 17:00, Tuesday

Research Session 10: Query Interfaces

Beacon A, 15:30 – 17:00, Tuesday Chair: Panos Ipeirotis

USHER: Improving Data Quality with Dynamic Forms

*Kuang Chen*¹, *Harr Chen*², *Neil Conway*¹, *Joseph M. Hellerstein*¹, *Tapan S. Parikh*¹; ¹*University of California at Berkeley, USA*; ²*MIT, USA*

Explaining Structured Queries in Natural Language

*Georgia Koutrika*¹, *Alkis Simitsis*², *Yannis E. Ioannidis*³; ¹*Stanford University, USA*; ²*HP, USA*; ³*University of Athens, Greece*

ScoreFinder: A Method for Collaborative Quality Inference on User-Generated Content (Short paper)

Yang Liao, Aaron Harwood, Kotagiri Ramamohanarao; University of Melbourne, Australia

IQ^{*p*}: **Incremental Query Construction, a Probabilistic Approach** (*Short paper*)

*Elena Demidova*¹, *Xuan Zhou*², *Wolfgang Nejdl*¹; ¹L3S Research Center, *Germany*; ²CSIRO, Australia

17:00 - 18:00 Poster Setup (Beacon AB)

18:00 - 20:00 All Posters Session and Reception (Beacon AB)

08:00-08:30 Continental Breakfast (Regency Foyer)

WEDNESDAY 08:30 - 10:00

Keynote 2

Beacon AB, 08:30 – 10:00, Wednesday Chair: Jayant Haritsa

DBMS: Lessons from the First 50 Years, Speculations for the Next 50

Jeffrey F. Naughton; University of Wisconsin-Madison, USA

10:00-10:30 Coffee Break (Regency Foyer)

WEDNESDAY 10:30-12:00

Demo Session 1B

Regency A, 10:30 – 12:00, Wednesday

Same as Demo Session 1A on previous pages

Research Session 11: Top-K Queries

Regency B, 10:30 – 12:00, Wednesday Chair: Ralf Schenkel

TASM: Top-k Approximate Subtree Matching

Nikolaus Augsten¹, Denilson Barbosa², Michael Böhlen¹, Themis Palpanas³; ¹Free University of Bozen-Bolzano, Italy; ²University of Alberta, Canada; ³University of Trento, Italy

Reverse Top-k Queries

*Akrivi Vlachou*¹, *Christos Doulkeridis*¹, *Yannis Kotidis*², *Kjetil Nørvåg*¹; ¹*NTNU*, *Norway*; ²*AUEB*, *Greece*

Top-K Aggregation Queries Over Large Networks (Short paper)

*Xifeng Yan*¹, *Bin He*², *Feida Zhu*³, *Jiawei Han*⁴; ¹*University of California at Santa Barbara, USA*; ²*IBM, USA*; ³*Singapore Management University, Singapore*; ⁴*University of Illinois at Urbana-Champaign, USA*

TopCells: Keyword-Based Search of Top-k Aggregated Documents in Text Cube (*Short paper*)

Bolin Ding, Bo Zhao, Cindy Xide Lin, Jiawei Han, Chengxiang Zhai; University of Illinois at Urbana-Champaign, USA

Research Session 12: Workflow and Workload Management

Regency C, 10:30 – 12:00, Wednesday Chair: Holger Schwarz

Optimizing ETL Workflows for Fault-Tolerance

Alkis Simitsis, Kevin Wilkinson, Umeshwar Dayal, Malu Castellanos; HP, USA

Q-Cop: Avoiding Bad Query Mixes to Minimize Client Timeouts Under Heavy Loads

Sean Tozer, Tim Brecht, Ashraf Aboulnaga; University of Waterloo, Canada

Admission Control Mechanisms for Continuous Queries in the Cloud (Short paper)

Lory Al Moakar¹, Panos K. Chrysanthis¹, Christine Chung², Shenoda Guirguis¹, Alexandros Labrinidis¹, Panayiotis Neophytou¹, Kirk Pruhs¹; ¹University of Pittsburgh, USA; ²Connecticut College, USA

Interaction-Aware Prediction of Business Intelligence Workload Completion Times (Short paper)

*Mumtaz Ahmad*¹, *Songyun Duan*², *Ashraf Aboulnaga*¹, *Shivnath Babu*²; ¹*University of Waterloo, Canada;* ²*Duke University, USA*

Research Session 13: Indexing and Hashing

Regency D, 10:30 – 12:00, Wednesday Chair: Paul Larson

Fast In-Memory XPath Search Using Compressed Indexes

Diego Arroyuelo¹, Francisco Claude², Sebastian Maneth³, Veli Mäkinen⁴, Gonzalo Navarro⁵, Kim Nguyễn³, Jouni Sirén⁴, Niko Välimäki⁴; ¹Yahoo!, Chile; ²University of Waterloo, Canada; ³NICTA, Australia; ⁴University of Helsinki, Finland; ⁵University of Chile, Chile

Hashing Tree-Structured Data: Methods and Applications

Shirish Tatikonda, Srinivasan Parthasarathy; Ohio State University, USA

Estimating the Compression Fraction of an Index Using Sampling (Short paper)

*Stratos Idreos*¹, *Raghav Kaushik*², *Vivek Narasayya*², *Ravishankar Ramamurthy*²; ¹*CWI*, *The Netherlands*; ²*Microsoft*, *USA*

The Hybrid-Layer Index: A Synergic Approach to Answering Top-*k* Queries in Arbitrary Subspaces (Short paper)

*Jun-Seok Heo*¹, *Junghoo Cho*², *Kyu-Young Whang*¹; ¹*KAIST, Korea*; ²*University of California at Los Angeles, USA*

Industry Session 2: Data, Data, and More Data

Regency EF, 10:30 – 12:00, Wednesday Chair: Surajit Chaudhuri

Data Cleansing as a Transient Service

Tanveer A. Faruquie, Hima Prasad K., L. Venkata Subramaniam, Mukesh Mohania, Girish Venkatachaliah, Shrinivas Kulkarni, Pramit Basu; IBM, India

XBRL Repository — An Industrial approach of Management of XBRL Documents

Zhen Hua Liu, Thomas Baby, Sriram Krishnamurthy, Ying Lu, Qin Yu, Anguel Novoselsky, Vikas Arora; Oracle, USA

Visualizing Large-Scale RDF Data Using Subsets, Summaries, and Sampling in Oracle

Seema Sundara, Medha Atre, Vladimir Kolovski, Souripriya Das, Zhe Wu, Eugene Inseok Chong, Jagannathan Srinivasan; Oracle, USA

12:00 - 14:00 Business Lunch (Beacon AB)

WEDNESDAY 14:00-15:30

Demo Session 2B

Regency A, 14:00 – 15:30, Wednesday

Same as Demo Session 2A on previous pages

Research Session 14: Scientific Data Mining

Regency B, 14:00 – 15:30, Wednesday Chair: Ambuj Singh

The Model-Summary Problem and a Solution for Trees

*Biswanath Panda*¹, *Mirek Riedewald*², *Daniel Fink*³; ¹*Google, USA*; ²*Northeastern University, USA*; ³*Cornell University, USA*

Efficient and Accurate Discovery of Patterns in Sequence Datasets

*Avrilia Floratou*¹, *Sandeep Tata*², *Jignesh M. Patel*¹; ¹*University of Wisconsin-Madison, USA*; ²*IBM, USA*

Mining Mutation Chains in Biological Sequences

*Chang Sheng*¹, *Wynne Hsu*¹, *Mong Li Lee*¹, *Joo Chuan Tong*², *See-Kiong* Ng²; ¹National University of Singapore, Singapore; ²Institute of Infocomm Research, Singapore

Research Session 15: Database Performance and Reliability

Regency C, 14:00 – 15:30, Wednesday Chair: Shahin Shayandeh

Exploring Power-Performance Tradeoffs in Database Systems

*Zichen Xu*¹, *Yi-Cheng Tu*¹, *Xiaorui Wang*²; ¹*University of South Florida, USA*; ²*University of Tennessee, USA*

Workload Driven Index Defragmentation

Vivek Narasayya, Manoj Syamala; Microsoft, USA

Impact of Disk Corruption on Open-Source DBMS

Sriram Subramanian, Yupu Zhang, Rajiv Vaidyanathan, Haryadi S. Gunawi, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, Jeffrey F. Naughton; University of Wisconsin-Madison, USA

Research Session 16: Spatial Databases

Regency D, 14:00 – 15:30, Wednesday Chair: Michael Böhlen

Locating Mapped Resources in Web 2.0

Dongxiang Zhang, Beng Chin Ooi, Anthony K.H. Tung; National University of Singapore, Singapore

Preference Queries in Large Multi-Cost Transportation Networks

*Kyriakos Mouratidis*¹, *Yimin Lin*¹, *Man Lung Yiu*²; ¹*Singapore Management University, Singapore*; ²*Hong Kong Polytechnic University, China*

Approximate String Search in Spatial Databases

*Bin Yao*¹, *Feifei Li*¹, *Marios Hadjieleftheriou*², *Kun Hou*¹; ¹*Florida State University, USA*; ²*AT&T Labs Research, USA*

Seminar 2

Regency EF, 14:00 – 15:30, Wednesday

Privacy in Data Publishing

*Johannes Gehrke*¹, *Daniel Kifer*², *Ashwin Machanavajjhala*³; ¹*Cornell University, USA*; ²*Pennsylvania State University, USA*; ³*Yahoo!, USA*

15:30-16:00 Coffee Break (Regency Foyer)

WEDNESDAY 16:00-17:30

Research Session 17: Sensor Networks

Regency B, 16:00 – 17:30, Wednesday Chair: Farnoush Banaei-Kashani

Global Iceberg Detection Over Distributed Data Streams

*Haiquan Zhao*¹, *Ashwin Lall*¹, *Mitsunori Ogihara*², *Jun Xu*¹; ¹*Georgia Institute of Technology, USA*; ²*University of Miami, USA*

Non-Dyadic Haar Wavelets for Streaming and Sensor Data

Chetan Gupta, Choudur Lakshminarayan, Song Wang, Abhay Mehta; HP, USA

Ratio Threshold Queries Over Distributed Data Sources (Short paper)

*Rajeev Gupta*¹, *Krithi Ramamritham*², *Mukesh Mohania*¹; ¹*IBM, India*; ²*IIT Bombay, India*

Probabilistic Top-*k* **Query Processing in Distributed Sensor Networks** *(Short paper)*

Mao Ye¹, Xingjie Liu¹, Wang-Chien Lee¹, Dik Lun Lee²; ¹Pennsylvania State University, USA; ²Hong Kong University of Science & Technology, China

Research Session 18: Query Optimization

Regency C, 16:00 – 17:30, Wednesday Chair: Jingren Zhou

Polynomial Heuristics for Query Optimization

Nicolas Bruno, César Galindo-Legaria, Milind Joshi; Microsoft, USA

Optimized Query Evaluation Using Cooperative Sorts

Yu Cao, Ramadhana Bramandia, Chee-Yong Chan, Kian-Lee Tan; National University of Singapore, Singapore

Generating Code for Holistic Query Evaluation

Konstantinos Krikellas, Stratis D. Viglas, Marcelo Cintra; University of Edinburgh, UK

Research Session 19: Graph Mining

Regency D, 16:00 – 17:30, Wednesday Chair: Wook-Shin Han

Finding Clusters in Subspaces of Very Large, Multi-Dimensional Datasets

*Robson L.F. Cordeiro*¹, *Agma J.M. Traina*¹, *Christos Faloutsos*², *Caetano Traina Jr.*¹; ¹Universidade de São Paulo, Brazil; ²Carnegie Mellon University, USA

Similarity Search on Supergraph Containment

*Haichuan Shang*¹, *Ke Zhu*¹, *Xuemin Lin*¹, *Ying Zhang*¹, *Ryutaro Ichise*²; ¹*University of New South Wales, Australia;* ²*National Institute of Informatics, Japan*

Finding Top-k Maximal Cliques in an Uncertain Graph (Short paper)

Zhaonian Zou, Jianzhong Li, Hong Gao, Shuo Zhang; Harbin Institute of Technology, China

Progressive Clustering of Networks Using Structure-Connected Order of Traversal (Short paper)

Dustin Bortner, Jiawei Han; University of Illinois at Urbana-Champaign, USA

Seminar 2 (cont.)

Regency EF, 16:00 – 17:30, Wednesday

Research Session 20: Parallel Processing

Beacon A, 16:00 - 17:30, Wednesday Chair: Mirek Riedewald

Osprey: Implementing MapReduce-Style Fault Tolerance in a Shared-Nothing Distributed Database

Christopher Yang, Christine Yen, Ceryen Tan, Samuel R. Madden; MIT, USA

FPGA Acceleration for the Frequent Item Problem

Jens Teubner, Rene Mueller, Gustavo Alonso; ETH Zürich, Switzerland

Estimating the Progress of MapReduce Pipelines (Short paper)

Kristi Morton, Abram Friesen, Magdalena Balazinska, Dan Grossman; University of Washington, USA

Scalable Distributed-Memory External Sorting (Short paper)

Mirko Rahn, Peter Sanders, Johannes Singler; Karlsruhe Institute of Technology, Germany

THURSDAY 08:30 - 10:00

Keynote 3

Beacon AB, 08:30 – 10:00, Thursday Chair: Gerhard Weikum

How New is the Cloud?

Donald Kossmann; ETH Zürich, Switzerland

10:00-10:30 Coffee Break (Regency Foyer)

THURSDAY 10:30 - 12:00

Research Session 21: Keyword Search

Regency B, 10:30 – 12:00, Thursday Chair: Zhen Liu

Supporting Top-K Keyword Search in XML Databases

Liang Jeff Chen, Yannis Papakonstantinou; University of California at San Diego, USA

Personalized Web Search with Location Preferences

*Kenneth Wai-Ting Leung*¹, *Dik Lun Lee*¹, *Wang-Chien Lee*²; ¹*Hong Kong University of Science & Technology, China*; ²*Pennsylvania State University, USA*

Fuzzy Matching of Web Queries to Structured Data (Short paper)

*Tao Cheng*¹, *Hady W. Lauw*², *Stelios Paparizos*²; ¹*University of Illinois at Urbana-Champaign, USA*; ²*Microsoft, USA*

Toward Industrial-Strength Keyword Search Systems Over Relational Data (Short paper)

Akanksha Baid, Ian Rae, AnHai Doan, Jeffrey F. Naughton; University of Wisconsin-Madison, USA

Research Session 22: Query Processing

Regency C, 10:30 – 12:00, Thursday Chair: Florian Waas

Efficient Processing of Substring Match Queries with Inverted q-Gram Indexes

*Younghoon Kim*¹, *Kyoung-Gu Woo*², *Hyoungmin Park*¹, *Kyuseok Shim*¹; ¹*Seoul National University, Korea;* ²*Samsung Electronics, Korea*

Progressive Result Generation for Multi-Criteria Decision Support Queries

Venkatesh Raghavan, Elke A. Rundensteiner; Worcester Polytechnic Institute, USA

Nb-GCLOCK: A Non-Blocking Buffer Management Based on the Generalized CLOCK

Makoto Yui¹, Jun Miyazaki², Shunsuke Uemura³, Hayato Yamana¹; ¹Waseda University, Japan; ²NAIST, Japan; ³Nara Sangyo University, Japan

Research Session 23: Web and Collaborative Applications

Regency D, 10:30 – 12:00, Thursday Chair: Zack Ives

Effective Automated Object Matching

Diego Zardetto¹, Monica Scannapieco¹, Tiziana Catarci²; ¹Istituto Nazionale di Statistica, Italy; ²Università di Roma "La Sapienza", Italy

Efficient Identification of Coupled Entities in Document Collections *(Short paper)*

*Nikos Sarkas*¹, *Albert Angel*¹, *Nick Koudas*¹, *Divesh Srivastava*²; ¹*University of Toronto, Canada;* ²*AT&T Labs Research, USA*

On Supporting Effective Web Extraction (Short paper)

*Wook-Shin Han*¹, *Wooseong Kwak*¹, *Hwanjo Yu*²; ¹*Kyungpook National University, Korea*; ²*POSTECH, Korea*

A Partial Persistent Data Structure to Support Consistency in Real-Time Collaborative Editing (Short paper)

Qinyi Wu¹, Calton Pu¹, João Eduardo Ferreiar²; ¹Georgia Institute of Technology, USA; ²Universidade de São Paulo, Brazil

Detecting Bursty Events in Collaborative Tagging Systems (Short paper) Junjie Yao, Bin Cui, Yuxin Huang, Yanhong Zhou; Peking University, China

Panel 1

Regency EF, 10:30 - 12:00, Thursday, Workshop

Cloudy Skies for Data Management

David Campbell¹, Brian Cooper², Dean Jacobs³, Ashok Joshi⁴, Volker Markl⁵, Srinivas Narayanan⁶; ¹Microsoft, USA; ²Yahoo!, USA; ³SAP, Germany; ⁴Oracle, USA; ⁵Technische Universität Berlin, Germany; ⁶Facebook, USA

Seminar 3

Beacon A, 10:30 - 12:00, Thursday

Representation, Composition and Application of Preferences in Databases

*Georgia Koutrika*¹, *Evaggelia Pitoura*², *Kostas Stefanidis*²; ¹*Stanford University, USA*; ²*University of Ioannina, Greece*

12:00 – 13:30 Lunch (on your own)

13:30 - 14:00 Poster Preparation (Beacon AB)

14:00-15:30 "All Posters" Session (Beacon AB)

15:30-16:00 Coffee Break (Regency Foyer)

THURSDAY 16:00 - 17:30

Research Session 24: Scientific Databases

Regency B, 16:00 – 17:30, Thursday Chair: Feifei Li

Credibility-Enhanced Curated Database: Improving the Value of Curated Databases

Qun Ni, Elisa Bertino; Purdue University, USA

UV-Diagram: A Voronoi Diagram for Uncertain Data

*Reynold Cheng*¹, *Xike Xie*¹, *Man Lung Yiu*², *Jinchuan Chen*³, *Liwen Sun*¹; ¹*University of Hong Kong, China;* ²*Hong Kong Polytechnic University, China;* ³*Renmin University of China, China*

Supporting Real-World Activities in Database Management Systems *(Short paper)*

Mohamed Y. Eltabakh, Walid G. Aref, Ahmed K. Elmagarmid, Yasin N. Silva, Mourad Ouzzani; Purdue University, USA

XML-Based Computation for Scientific Workflows (Short paper)

Daniel Zinn¹, Shawn Bowers², Bertram Ludäscher¹; ¹University of California at Davis, USA; ²Gonzaga University, USA

Research Session 25: Tree Queries and Semi-Structured Databases

Regency C, 16:00 – 17:30, Thursday Chair: Maurice van Keulen

ViewJoin: Efficient View-Based Evaluation of Tree Pattern Queries

Ding Chen, Chee-Yong Chan; National University of Singapore, Singapore

FlexPref: A Framework for Extensible Preference Evaluation in Database Systems

Justin J. Levandoski, Mohamed F. Mokbel, Mohamed E. Khalefa; University of Minnesota, USA

Optimal Tree Node Ordering for Child/Descendant Navigations *(Short paper)*

*Atsuyuki Morishima*¹, *Keishi Tajima*², *Masateru Tadaishi*¹; ¹*University of Tsukuba, Japan*; ²*Kyoto University, Japan*

XMorph: A Shape-Polymorphic, Domain-Specific XML Data Transformation Language (Short paper)

*Curtis Dyreson*¹, *Sourav Bhowmick*², *Aswani Rao Jannu*¹, *Kirankanth Mallampalli*¹, *Shuohao Zhang*³; ¹Utah State University, USA; ²Nanyang Technological University, Singapore; ³Marvel, USA

Research Session 26: Query Ranking and Database Testing

Regency D, 16:00 – 17:30, Thursday Chair: Cesar Galindo-Legaria

Surrogate Ranking for Very Expensive Similarity Queries

*Fei Xu*¹, *Ravi Jampani*¹, *Mingxi Wu*², *Chris Jermaine*¹, *Tamer Kahveci*¹; ¹*University of Florida, USA*; ²*Oracle, USA*

Semantic Ranking and Result Visualization for Life Sciences Publications

Julia Stoyanovich, William Mee, Kenneth A. Ross; Columbia University, USA

Ranked Queries Over Sources with Boolean Query Interfaces without Ranking Support (Short paper)

*Vagelis Hristidis*¹, *Yuheng Hu*¹, *Panagiotis G. Ipeirotis*²; ¹*Florida International University, USA*; ²*New York University, USA*

X-Data: Generating Test Data for Killing SQL Mutants (Short paper)

Bhanu Pratap Gupta, Devang Vira, S. Sudarshan; IIT Bombay, India

Seminar 4

Regency EF, 16:00 – 17:30, Thursday

Database as a Service (DBaaS)

*Wolfgang Lehner*¹, *Kai-Uwe Sattler*²; ¹*Dresden University of Technology, Germany*; ²*Ilmenau University of Technology, Germany*

Panel 2

Beacon A, 16:00 - 17:30, Thursday, Workshop

Database Architecture (R)evolution: New Hardware vs. New Software

*Stavros Harizopoulos*¹, *Tassos Argyros*², *Peter A. Boncz*³, *Dan Dietterich*⁴, *Samuel R. Madden*⁵, *Florian M. Waas*⁶; ¹*HP*, *USA*; ²*Aster Data, USA*; ³*CWI*, *The Netherlands*; ⁴*Netezza, USA*; ⁵*MIT, USA*; ⁶*Greenplum, USA*

18:00 - 19:00 Pre-Banquet Presentation (Regency Foyer)

Double Blind Revue Jazz Band

The presentation will cover string algorithms, horn clauses, foreign keys, data bass management, and drum scheduling.

THURSDAY 19:00 - 21:00

Banquet

Beacon AB, 19:00 - 21:00, Thursday

Banquet Speaker: *Gio Wiederhold, Stanford University, USA* **"From Crossing Chasms to Climbing into Clouds"**
FRIDAY 08:30 - 10:00

Research Session 27: Social Networks and Similarity Queries

Regency B, 08:30 – 10:00, Friday Chair: K. Selcuk Candan

Discovery-Driven Graph Summarization

*Ning Zhang*¹, *Yuanyuan Tian*², *Jignesh M. Patel*¹; ¹*University of Wisconsin-Madison, USA*; ²*IBM, USA*

The Similarity Join Database Operator

*Yasin N. Silva*¹, *Walid G. Aref*¹, *Mohamed H. Ali*²; ¹*Purdue University, USA*; ²*Microsoft, USA*

Anonymizing Weighted Social Network Graphs (Short paper)

Sudipto Das, Ömer Eğecioğlu, Amr El Abbadi; University of California at Santa Barbara, USA

Efficient Similarity Matching of Time Series Cliques with Natural Relations *(Short paper)*

*Zhe Zhao*¹, *Bin Cui*¹, *Wee Hyong Tok*², *Jiakui Zhao*³; ¹*Peking University, China*; ²*Microsoft, China*; ³*China Electric Power Research Institute, China*

Research Session 28: Stream Processing

Regency C, 08:30 – 10:00, Friday Chair: Alexandros Labrinidis

Continuous Query Evaluation Over Distributed Sensor Networks

Oana Jurca, Sebastian Michel, Alexandre Herrmann, Karl Aberer; EPFL, Switzerland

Space-Efficient Online Approximation of Time Series Data: Streams, Amnesia, and Out-of-Order

Sorabh Gandhi, Luca Foschini, Subhash Suri; University of California at Santa Barbara, USA

Approximation Trade-Offs in Markovian Stream Processing: An Empirical Study (Short paper)

*Julie Letchner*¹, *Christopher Ré*², *Magdalena Balazinska*¹, *Matthai Philipose*³; ¹*University of Washington, USA*; ²*University of Wisconsin-Madison, USA*; ³*Intel, USA*

FENCE: Continuous Access Control Enforcement in Dynamic Data Stream Environments (*Short paper*)

*Rimma V. Nehme*¹, *Hyo-Sang Lim*², *Elisa Bertino*²; ¹*Microsoft, USA*; ²*Purdue University, USA*

ICDE 2010

Industry Session 3: Query Optimization

Regency D, 08:30 - 10:00, Friday Chair: Vinavak Borkar

Incorporating Partitioning and Parallel Plans into the SCOPE Optimizer Jingren Zhou, Per-Ake Larson, Ronnie Chaiken; Microsoft, USA

Rule Profiling for Query Optimizers and Their Implications

Surajit Chaudhuri, Leo Giakoumakis, Vivek Narasayya, Ravishankar Ramamurthy; Microsoft, USA

Data Desensitization of Customer Data for Use in Optimizer Performance Experiments

*Malu Castellanos*¹, *Bin Zhang*¹, *Ivo Jimenez*¹, *Perla Ruiz*², *Miguel* Durazo², Umeshwar Dayal¹, Lily Jow¹; ¹HP, USA; ²University of Sonora, Mexico

Seminar 5

Regency EF, 08:30 - 10:00, Friday

Techniques for Efficiently Searching in Spatial, Temporal, Spatio-Temporal, and Multimedia Databases

Hans-Peter Kriegel, Peer Kröger, Matthias Renz; LMU München, Germany

PhD Workshop 1: Applications, Uncertainty, and Privacy

Beacon A, 08:30 - 10:00, Friday, Workshop

Maximizing Visibility of Objects

Muhammed Miah; University of Texas at Arlington, USA

Improving Product Search with Economic Theory

Beibei Li; New York University, USA

Toward Large Scale Data-Aware Search: Ranking, Indexing, Resolution and Beyond

Tao Cheng; University of Illinois at Urbana-Champaign, USA

Graphical Models for Dependencies and Queries in Uncertain Data

Ruiwen Chen; University of Ottawa, Canada

Privacy-Preserving Data Publishing Ruilin Liu; Stevens Institute of Technology, USA

10:00 – 10:30 Coffee Break (Regency Foyer)

FRIDAY 10:30-12:00

Research Session 29: Publishing Privacy

Regency B, 10:30 – 12:00, Friday Chair: Brian Cooper

A Privacy-Preserving Approach to Policy-Based Content Dissemination

Ning Shang, Mohamed Nabeel, Federica Paci, Elisa Bertino; Purdue University, USA

Global Privacy Guarantee in Serial Data Publishing (Short paper)

*Raymond Chi-Wing Wong*¹, *Ada Wai-Chee Fu*², *Jia Liu*², *Ke Wang*³, *Yabo Xu*⁴; ¹*Hong Kong University of Science & Technology, China*; ²*Chinese University of Hong Kong, China*; ³*Simon Fraser University, Canada*; ⁴*Sun Yat-sen University, China*

XColor: Protecting General Proximity Privacy (Short paper)

Ting Wang, Ling Liu; Georgia Institute of Technology, USA

Correlation Hiding by Independence Masking (Short paper)

Yufei Tao¹, Jian Pei², Jiexing Li¹, Xiaokui Xiao³, Ke Yi⁴, Zhengzheng Xing²; ¹Chinese University of Hong Kong, China; ²Simon Fraser University, Canada; ³Nanyang Technological University, Singapore; ⁴Hong Kong University of Science & Technology, China

Research Session 30: Data Clouds

Regency C, 10:30 – 12:00, Friday Chair: Walid Aref

Monitoring Continuous State Violation in Datacenters: Exploring the Time Dimension

Shicong Meng, Ting Wang, Ling Liu; Georgia Institute of Technology, USA

Cost-Efficient and Differentiated Data Availability Guarantees in Data Clouds (Short paper)

Nicolas Bonvin, Thanasis G. Papaioannou, Karl Aberer; EPFL, Switzerland

Intensional Associations in Dataspaces (Short paper) Marcos Antonio Vaz Salles¹, Jens Dittrich², Lukas Blunschi³; ¹Cornell University, USA; ²Saarland University, Germany; ³ETH Zürich, Switzerland

A Tuple Space for Social Networking on Mobile Phones (Short paper) Emre Sarigöl, Oriana Riva, Gustavo Alonso; ETH Zürich, Switzerland

Overlapping Community Search for Social Networks (Short paper) Arnau Padrol-Sureda, Guillem Perarnau-Llobet, Julian Pfeifle, Victor Muntés-Muleor; Universitat Politècnica de Catalunya, Spain

Seminar 5 (cont.)

Regency EF, 10:30 – 12:00, Friday

PhD Workshop 2: Storage, Indexing, and Search

Beacon A, 10:30 – 12:00, Friday, Workshop

Flash-Enabled Database Storage

Ioannis Koltsidas; University of Edinburgh, UK

A Database Server for Next-Generation Scientific Data Management

Mohamed Y. Eltabakh; Purdue University, USA

CareDB: A Context and Preference-Aware Location-Based Database System

Justin J. Levandoski; University of Minnesota, USA

Graph Indexing for Reachability Queries

Hilmi Yıldırım; Rensselaer Polytechnic Institute, USA

Evaluating Path Queries Over Route Collections

Panagiotis Bouros; National Technical University of Athens, Greece

12:00 - 13:30 Lunch (on your own)

FRIDAY 13:30 - 15:00

DESWeb Keynote and Panel Regency B, 13:30 – 15:00, Friday

Keynote Speaker: Howard Ho (IBM Almaden)

Keynote Title: Integrating Linked Open Data Sources about Financial Companies: Experience and Lessons Learned

Panel: TBA

NTII Opening Remarks and Keynote 1

Regency EF, 13:30 – 15:00, Friday

Keynote Speaker: Dan Wolfson (IBM)

Keynote Title: Business Information Management and Controls: Lessons from the Current Financial Crisis

PhD Workshop 3: Data Mining

Beacon A, 13:30 – 15:00, Friday, Workshop

Advances in Constrained Clustering

ZiJie Qi; University of California at Davis, USA

Towards a Task-Based Search and Recommender Systems

Gabriele Tolomei; Università Ca' Foscari Venezia, Italy

On Dynamic Data Clustering and Visualization Using Swarm Intelligence

Esin Saka; University of Louisville, USA

Fast Algorithms for Time Series Mining

Lei Li; Carnegie Mellon University, USA

15:00-15:30 Coffee Break (Regency Foyer)

FRIDAY 15:30-17:00

DESWeb Session 1: Data Search

Regency B, 15:30 - 17:00, Friday, Workshop

Semantic Flooding: Search Over Semantic Links

*Fausto Giunchiglia*¹, *Uladzimir Kharkevich*¹, *Alethia Hume*¹, *Piyatat Chatvorawit*²; ¹*University of Trento, Italy*; ²*Asian Institute of Technology, Thailand*

An Ontology-Based Retrieval System Using Semantic Indexing

*Soner Kara*¹, *Özgür Alan*¹, *Orkunt Sabuncu*¹, *Samet Akpınar*², *Nihan K. Çiçekli*², *Ferda N. Alpaslan*²; ¹*Orbim Corp., Turkey*; ²*METU, Turkey*

Keyword Based Search Over Semantic Data in Polynomial Time

*Paolo Cappellari*¹, *Roberto De Virgilio*², *Antonio Maccioni*², *Michele Miscione*²; ¹*University of Alberta, Canada;* ²*Università Roma Tre, Italy*

NTII Session 1: New Architectures and Models

Regency EF, 15:30 – 17:00, Friday, Workshop

A First Step Towards Integration Independence

*Laura M. Haas*¹, *Renée J. Miller*², *Donald Kossmann*³, *Martin Hentschel*³; ¹*IBM, USA;* ²*University of Toronto, Canada;* ³*ETH Zürich, Switzerland*

Streaming Data Integration: Challenges and Opportunities *Nesime Tatbul; ETH Zürich, Switzerland*

Partitioning Real-Time ETL Workflows

Alkis Simitsis, Chetan Gupta, Song Wang, Umeshwar Dayal; HP, USA

FRIDAY 17:00-18:00

DESWeb Session 2: Techniques for Mappings

Regency B, 17:00 - 18:00, Friday, Workshop

Ontology Alignment Argumentation with Mutual Dependency Between Arguments and Mappings

Paulo Maio, Nuno Silva; Politécnico do Porto, Portugal

Summarizing Ontology-Based Schemas in PDMS

*Carlos Eduardo Pires*¹, *Paulo Sousa*², *Zoubida Kedad*³, *Ana Carolina Salgado*²; ¹*UFCG, Brazil*; ²*UFPE, Brazil*; ³*UVSQ, France*

A Framework for Automatic Schema Mapping Verification Through Reasoning

Paolo Cappellari¹, Denilson Barbosa¹, Paolo Atzeni²; ¹University of Alberta, Canada; ²Università Roma Tre, Italy

08:00-08:30 Continental Breakfast (Regency Foyer)

SATURDAY 08:30 - 10:00

DESWeb Session 3: Data Management

Regency B, 09:10 – 10:00, Saturday, Workshop

Extensions to the Pig Data Processing Platform for Scalable RDF Data Processing Using Hadoop

Yusuke Tanimura, Akiyoshi Matono, Steven Lynden, Isao Kojima; AIST, Japan

Optimized Data Access for Efficient Execution of Semantic Services

Thorsten Möller, Heiko Schuldt; University of Basel, Switzerland

M3SN Keynote and Session 1

Regency C, 08:30 - 10:00, Saturday, Workshop

Keynote Speaker: TBA

Privometer: Privacy Protection in Social Networks

Nilothpal Talukder, Mourad Ouzzani, Ahmed K. Elmagarmid, Hazem Elmeleegy, Mohamed Yakout; Purdue University, USA

NTII Keynote 2

Regency EF, 08:30 - 10:00, Saturday, Workshop

Keynote Speaker: *Craig Knoblock (USC)*

Keynote Title: Interactively Building Geospatial Mashups

10:00-10:30 Coffee Break (Regency Foyer)

SATURDAY 10:30-12:00

DESWeb Session 4: Entity Analysis Regency B, 10:30 - 12:00, Saturday, Workshop

The Entity Name System: Enabling the Web of Entities

Heiko Stoermer, Themis Palpanas, George Giannakopoulos; University of Trento, Italy

Towards Better Entity Resolution Techniques for Web Document Collections

Surender Reddy Yerva, Zoltán Miklós, Karl Aberer; EPFL, Switzerland

Processing Online News Streams for Large-Scale Semantic Analysis

*Miloš Krstajić*¹, *Florian Mansmann*¹, *Andreas Stoffel*¹, *Martin Atkinson*², *Daniel A. Keim*¹; ¹*University of Konstanz, Germany*; ²*EC Joint Research Centre, Italy*

DIVERSUM: Towards Diversified Summarisation of Entities in Knowledge Graphs

*Marcin Sydow*¹, *Mariusz Pikuła*¹, *Ralf Schenkel*²; ¹*Polish-Japanese Institute of Information Technology, Poland*; ²*Max-Planck Institute for Informatics, Germany*

M3SN Session 2

Regency C, 10:30 – 12:00, Saturday, Workshop

On the Influence of Social Factors on Team Recommendations

Michele Brocco, Georg Groh, Christian Kern; Technische Universität München, Germany

Towards Discovery of Eras in Social Networks

*Michele Berlingerio*¹, *Michele Coscia*¹, *Fosca Giannotti*¹, *Anna Monreale*¹, *Dino Pedreschi*²; ¹*CNR-ISTI, Italy*; ²*University of Pisa, Italy*

Mining and Representing Recommendations in Actively Evolving Recommender Systems

Ira Assent; Aalborg University, Denmark

NTII Session 2: Applications

Regency EF, 10:30 – 12:00, Saturday, Workshop

Towards Best-Effort Merge of Taxonomically Organized Data

David Thau¹, Shawn Bowers², Bertram Ludäscher¹; ¹University of California at Davis, USA; ²Gonzaga University, USA

BI-Style Relation Discovery Among Entities in Text

Wojciech M. Barczyński, Falk Brauer, Adrian Mocan, Marcus Schramm, Jan Froemberg; SAP, Germany

Profiling Linked Open Data with ProLOD

Christoph Böhm, Felix Naumann, Ziawasch Abedjan, Dandy Fenz, Toni Grütze, Daniel Hefenbrock, Matthias Pohl, David Sonnabend; HPI, Germany

Coordination of Data in Heterogenous Domains

Michael Lawrence, Rachel Pottinger, Sheryl Staub-French; University of British Columbia, Canada

SATURDAY 13:30-15:00

NTII Session 3 and Wrap-Up: New Primitives and Techniques

Regency EF, 13:30 - 15:00, Saturday, Workshop

Duplicate Detection in Probabilistic Data

*Fabian Panse*¹, *Maurice van Keulen*², *Ander de Keijzer*², *Norbert Ritter*¹; ¹*University of Hamburg, Germany;* ²*University of Twente, The Netherlands*

Complement Union for Data Integration

*Jens Bleiholder*¹, *Sascha Szott*², *Melanie Herschel*³, *Felix Naumann*¹; ¹*HPI, Germany*; ²*ZIB, Germany*; ³*Universität Tübingen, Germany*

Midas for Government: Integration of Government Spending Data on Hadoop

Antonio Sala¹, Calvin Lin², Howard Ho³; ¹University of Modena & Reggio Emilia, Italy; ²University of California at Berkeley, USA; ³IBM, USA

ICDE 2010 Keynotes

Keynote 1

Beacon AB, 08:30 - 10:00 Tuesday 2nd March

Large Scale Data Warehousing: Trends and Observations

*Richard Winter*¹, *Pekka Kostamaa*²; ¹*Winter Corporation, USA*; ²*Teradata, USA*

How large are data warehouses? How fast are they growing? How big are they going to get? What is driving their growth? Why is all this data of value in commercial enterprises? What can we say about how these large data warehouses are being used? What are some key challenges ahead? In this talk, Richard Winter will share his views and observations concerning these questions and others, based on more than three decades of involvement with commercial data warehouses and their preferences.





Richard Winter is an industry expert in large scale data management technology, architecture and implementation with over twenty-five years of experience. As President of WINTERCORP (www.wintercorp.com), a consulting firm in Cambridge MA, he advises executives on their strategies and critical projects, focusing on scalability, performance and availability in data warehousing. Mr. Winter is a frequent author and speaker; and, teaches seminars on scalability management, including the architecture and selection of data warehouse platforms. From 1995–2005, Mr. Winter founded and directed the Winter Large Scale Database Research Project, which measured the scale of the largest publicly acknowledged databases in the world and published its findings on the internet and in a series of research reports.

Pekka Kostamaa is the Director of Advanced Development and Enterprise Architecture for Teradata R&D. His group is responsible for research of advanced topics for features to be added to the Teradata Warehouse platform in the 3–5 year time frame. The group also produces Intellectual Property and maintains academic and university relationships for the company.

Pekka has been involved with the Teradata database for over 20 years in several different development functions. Among others, he has been a performance analyst, a software developer, a development architect, and a department manager for the performance and optimizer groups. Before Teradata and during a few year absence from Teradata, Pekka joined some start-up companies, and did consulting in the data warehouse and database areas. He has several publications, holds six patents with about 20 pending, and is a member of the UCLA Computer Science Advisory Board.

Keynote 2

Beacon AB, 08:30 - 10:00 Wednesday 3rd March

DBMS: Lessons from the First 50 Years, Speculations for the Next 50

Jeffrey F. Naughton; University of Wisconsin-Madison, USA

Some of the major themes in DBMS research appeared in the computer science literature as early as 50 years ago. The community has had a very productive time over the past 50 years exploring these themes, in the process contributing to a major software industry and creating a large and vibrant research community. I will give a subjective and probably highly biased view of these themes and why they have been so persistent, and speculate on how they might continue to persist in the future. While we will probably continue to be

productive over the next 50 years as well, there are reasons for concern going forward. I will close with some speculation on what we might do to deal with this.



Jeffrey F. Naughton earned a bachelor's degree in Mathematics from the University of Wisconsin-Madison, and a PhD in Computer Science from Stanford University. He served as a faculty member in the Computer Science Department of Princeton University before moving to the University of Wisconsin-Madison, where he is currently Professor of Computer Science. He received the National Science Foundation's Presidential Young Investigator Award in 1991, received the Vilas Associate Award for Excellence in Research in 2000, was named an ACM Fellow in 2002, and was a member of the Gamma project that received the 2008 ACM Software Systems Award. He has served as a consultant for companies including Greenplum and Teradata/NCR, and is currently a consultant at the Microsoft Jim Gray Systems Lab.

Keynote 3

Beacon AB, 08:30 - 10:00 Thursday 4th March

How New is the Cloud?

Donald Kossmann; ETH Zürich, Switzerland

Cloud computing has been identified by Gartner as one of the ten most disruptive technologies for the next decade. It has made many promises and the first products have appeared on the market place and are rapidly gaining adoption. Time to step back a bit.

This talk first gives an overview of the promises made by cloud computing. Which promises really matter? Which promises were only made because they could be fulfilled? And, which promises were only made because they could not be validated? Second, this talk discusses the fundamental limitations, light-housed by the CAP theorem. How bad is it really? Third, this talk discusses alternative architectures for data management in the cloud. What works? What is new? Fourth, this talk addresses the changes application programmers will face. What is realistic? Finally, this talk reports on three years of running a cloud computing start-up. What do users like, what do investors like? What do I like?



Donald Kossmann is a professor for Computer Science at ETH Zurich (Switzerland) and the CEO of 28msec Inc. He received his MS in 1991 from the University of Karlsruhe and completed his PhD in 1995 at the Technical University of Aachen. After that, he held positions at the University of Maryland, the IBM Almaden Research Center, the University of Passau, the Technical University of Munich, and the University of Heidelberg. At ETH Zurich and 28msec, he develops new technologies at the intersection of database systems, web technologies, and distributed systems. Before joining ETH and 28msec, Donald Kossmann was a co-founder of i-TV-T AG (1998, still in business) and XQRL Inc. which was founded in 2002 and acquired by BEA in the same year.

ICDE 2010 BANQUET

The conference banquet will be held in Beacon AB at 19:00 - 21:00 on Thursday 4th March. A pre-banquet jazz presentation by the band Double Blind Revue will be held in the Regency Foyer from 18:00 - 19:00.

Pre-Banquet Presentation

Regency Foyer, 18:00 - 19:00 Thursday March 4th

Double Blind Revue Jazz Band

The presentation will cover string algorithms, horn clauses, foreign keys, data bass management, and drum scheduling.

Banquet Speaker

From Crossing Chasms to Climbing into Clouds

Gio Wiederhold, Stanford University, USA

A prime motivation in 1984 for initiating the IEEE Data Engineering Conferences was to bridge the gap from theory to practice. I'll briefly summarize at a high level phases in system approaches from those early days to today: Relational databases replacing established data management systems, moving on to the use of clouds where the storage responsibilities are fully delegated. Parallel to those developments access for processing changed from batchprocessing to timesharing, with competition from free-standing systems, and soon after combinations as client-server approaches.

Many of those changes are motivated by economic considerations: the cost of hardware, software, and people, and those costs diverged. I believe that computer science community is doing itself a disservice by not understanding economic factors better. Today we leave it to others, businesses, economists and promoters to value and sell our products, and then complain about the results in terms of opportunities. Engineering remains at the forefront of bringing innovations to practice. To remain successful the computing community has to understand better than it does now what it contributes and how it will be valued.



Gio Wiederhold is now an Emeritus Professor of Computer Science, Electrical Engineering, and Medicine at Stanford University, continuing part-time with courses on 'Business on the Internet' and 'Software Economics'. Gio Wiederhold was born in Italy in 1936, received a degree in Aeronautical Engineering in Holland in 1957 and a PhD in Medical Information Science from the University of California at San Francisco in 1976.

Phases of his life included 16 years in industry, 22 years fulltime academia, and 10 years government service, with crossinteractions throughout. Research topics addressed combustion analysis, compilers, timesharing, database technologies, knowledge-based integration of information, an algebra over ontologies, access to simulations, privacy protection, composition of software, and semantic complexity of web.

Gio Wiederhold has been elected fellow of the ACMI, the IEEE, and the ACM. He has been an editor and editor-in-chief of several IEEE and ACM publications. Gio's web page is at http://infolab.stanford.edu/people/gio.html.

ICDE 2010 Seminars

SEMINAR 1

Regency EF, 13:30 - 17:00 Tuesday March 2nd 2010

Anonymized Data: Generation, Models, Usage

Graham Cormode, Divesh Srivastava; AT&T Labs Research, USA

Data anonymization techniques enable publication of detailed information, which permits ad hoc queries and analyses, while guaranteeing the privacy of sensitive information in the data against a variety of attacks. In this tutorial, we aim to present a *unified* framework of data anonymization techniques, viewed through the lens of data uncertainty. Essentially, anonymized data describes a set of possible worlds that include the original data. We show that anonymization approaches generate different working models of uncertain data, and that the privacy guarantees offered by *k*-anonymization and *l*-diversity can be naturally understood in terms of the sets of possible worlds that correspond to the anonymized data. We show the in query evaluation over uncertain databases can hence be used for answering ad hoc queries over anonymized data. We identify new research problems for both the Data Anonymization and the Uncertain Data communities.

SEMINAR 2

Regency EF, 14:00 - 17:30 Wednesday March 3rd 2010

Privacy in Data Publishing

*Johannes Gehrke*¹, Daniel Kifer², Ashwin Machanavajjhala³; ¹Cornell University, USA; ²Pennsylvania State University, USA; ³Yahoo!, USA

This tutorial gives an overview of techniques for releasing data about individuals while preserving privacy.

SEMINAR 3

Beacon A, 10:30 - 12:00 Thursday March 4th 2010

Representation, Composition and Application of Preferences in Databases

*Georgia Koutrika*¹, *Evaggelia Pitoura*², *Kostas Stefanidis*²; ¹*Stanford University, USA*; ²*University of Ioannina, Greece*

This tutorial provides an overview of the key research results in the area of user preferences from a database perspective. The objective is to survey in a systematic and holistic way a number of approaches for preference representation and composition, querying with preferences and preference learning. Open research problems are also presented.

SEMINAR 4

Regency EF, 16:00 - 17:30 Thursday March 4th 2010

Database as a Service (DBaaS)

*Wolfgang Lehner*¹, *Kai-Uwe Sattler*²; ¹*Dresden University of Technology, Germany*; ²*Ilmenau University of Technology, Germany*

Modern Web or "Eternal-Beta" applications necessitate a flexible and easy-to-use data management platform that allows the evolutionary development of databases and applications. The classical approach of relational database systems following strictly the ACID properties has to be extended by an extensible and easy-to-use persistency layer with specialized DB features. Using the underlying concept of Software as a Service (SaaS) also enables an economic advantage based on the "economy of the scale", where application and system environments only need to be provided once but can be used by thousands of users. Within this tutorial, we are looking at the current state-of-the-art from different perspectives. We outline foundations and techniques to build database services based on the SaaS-paradigm. We discuss requirements from a programming perspective, show different dimensions in the context of consistency and reliability, and also describe different non-functional properties under the umbrella of Service-Level agreements (SLA).

SEMINAR 5

Regency EF, 08:30 - 12:00 Friday March 5th 2010

Techniques for Efficiently Searching in Spatial, Temporal, Spatio-Temporal, and Multimedia Databases

Hans-Peter Kriegel, Peer Kröger, Matthias Renz; LMU München, Germany

This tutorial provides a comprehensive and comparative overview of general techniques to efficiently support similarity queries in spatial, temporal, spatio-temporal, and multimedia databases. In particular, it identifies the most generic query types and discusses general algorithmic methods to answer such queries efficiently. In addition, the tutorial sketches important applications of the introduced methods, and presents sample implementations of the general approaches within each of the aforementioned database types. The intended audience of this tutorial ranges from novice researchers to advanced experts as well as practitioners from any application domain dealing with spatial, temporal, spatio-temporal, and/or multimedia data.

Research Session 1: KNN Queries

Regency B, 10:30 – 12:00, Tuesday Chair: Julia Stoyanovich

K Nearest Neighbor Queries and KNN-Joins in Large Relational Databases (Almost) for Free

Bin Yao, Feifei Li, Piyush Kumar; Florida State University, USA

Finding the *k* nearest neighbors (*k*NN) of a query point, or a set of query points (*k*NN-Join) are fundamental problems in many application domains. Many previous efforts to solve these problems focused on spatial databases or stand-alone systems, where changes to the database engine may be required, which may limit their application on large data sets that are stored in a relational database management system. Furthermore, these methods may not automatically optimize kNN queries or kNN-Joins when additional query conditions are specified. In this work, we study both the kNN query and the kNN-Join in a relational database, possibly augmented with additional query conditions. We search for relational algorithms that require no changes to the database engine. The straightforward solution uses the user-defined-function (UDF) that a query optimizer cannot optimize. We design algorithms that could be implemented by SQL operators without changes to the database engine, hence enabling the query optimizer to understand and generate the "best" query plan. Using only a small constant number of random shifts for databases in any fixed dimension, our approach guarantees to find the approximate kNN with only logarithmic number of page accesses in expectation with a constant approximation ratio and it could be extended to find the exact kNN efficiently in any fixed dimension. Our design paradigm easily supports the kNN-Join and updates. Extensive experiments on large, real and synthetic, data sets confirm the efficiency and practicality of our approach.

Quantile-Based KNN Over Multi-Valued Objects

Wenjie Zhang, Xuemin Lin, Muhammad Aamir Cheema, Ying Zhang, Wei Wang; University of New South Wales, Australia

K Nearest Neighbor search has many applications including data mining, multi-media, image processing, and monitoring moving objects. In this paper, we study the problem of *K*NN over multi-valued objects. We aim to provide effective and efficient techniques to identify KNN sensitive to relative distributions of objects. We propose to use quantiles to summarize relative-distribution-sensitive *K* nearest neighbors. Given a query *Q* and a quantile $\phi \in (0, 1]$, we firstly study the problem of efficiently computing *K* nearest objects based on a ϕ -quantile distance (e.g. median distance) from each object to *Q*. The second problem is to retrieve the *K* nearest objects to *Q* based on overall distances in the "best population" (with a given size specified by ϕ -quantile) for each object. While the first problem can be solved in polynomial time, we show that the 2nd problem is NP-hard. A set of efficient, novel algorithms have been proposed to give an exact solution for the first problem and an approximate solution for the second problem with the approximation ratio 2. Extensive experiment demonstrates that our techniques are very efficient and effective.

Efficient Rank Based KNN Query Processing Over Uncertain Data

Ying Zhang, Xuemin Lin, Gaoping Zhu, Wenjie Zhang, Qianlu Lin; University of New South Wales, Australia

Uncertain data are inherent in many applications such as environmental surveillance and quantitative economics research. As an important problem in many applications, KNN query has been extensively investigated in the literature. In this paper, we study the problem of processing rank based KNN query against uncertain data. Besides applying the *expected rank* semantic to compute KNN, we also introduce the *median rank* which is less sensitive to the outliers. We show both ranking methods satisfy nice top-*k* properties such as *exact-k*, *containment*, *unique ranking*, *value invariance*, *stability* and *fairfulness*. For

given query q, IO and CPU efficient algorithms are proposed in the paper to compute KNN based on expected (median) ranks of the uncertain objects. To tackle the correlations of the uncertain objects and high IO cost caused by large number of instances of the uncertain objects, randomized algorithms are proposed to approximately compute KNN with theoretical guarantees. Comprehensive experiments are conducted on both real and synthetic data to demonstrate the efficiency of our techniques.

Research Session 2: Distributed Data

Regency C, 10:30 – 12:00, Tuesday Chair: Hank Korth

Reliable Storage and Querying for Collaborative Data Sharing Systems

Nicholas E. Taylor, Zachary G. Ives; University of Pennsylvania, USA

The sciences, business confederations, and medicine urgently need infrastructure for sharing data and updates among collaborators' constantly changing, heterogeneous databases. The ORCHESTRA system addresses these needs by providing data transformation and exchange capabilities across DBMSs, combined with archived storage of all database versions. ORCHESTRA adopts a peer-to-peer architecture in which individual collaborators contribute data and compute resources, but where there may be no dedicated server or compute cluster.

We study how to take the combined resources of ORCHESTRA's autonomous nodes, as well as PCs from "cloud" services such as Amazon EC2, and provide reliable, *cooperative* storage and query processing capabilities. We guarantee reliability and correctness as in distributed or cloud DBMSs, while also supporting cross-domain deployments, replication, and transparent failover, as provided by peer-to-peer systems. Our storage and query subsystem supports dozens to hundreds of nodes across different domains, possibly including nodes on cloud services.

Our contributions include (1) a modified data partitioning substrate that combines cluster and peer-to-peer techniques, (2) an efficient implementation of replicated, reliable, versioned storage of relational data, (3) new query processing and indexing techniques over this storage layer, and (4) a mechanism for incrementally recomputing query results that ensures correct, complete, and duplicate-free results in the event of node failure during query execution. We experimentally validate query processing performance, failure detection methods, and the performance benefits of incremental recovery in a prototype implementation.

Strongly Consistent Replication for a Bargain

Konstantinos Krikellas¹, Sameh Elnikety², Zografoula Vagena³, Orion Hodson²; ¹University of Edinburgh, UK; ²Microsoft, UK; ³Concentra Consulting Ltd., UK

Strong consistency is an important correctness property for replicated databases. It ensures that each transaction accesses the latest committed database state as provided in centralized databases. Achieving strong consistency in replicated databases is a major performance challenge and is typically not provided, exposing inconsistent data to client applications. We propose two scalable techniques that exploit lazy update propagation and workload information to guarantee strong consistency by delaying transaction start. We implement a prototype replicated database system and incorporate the proposed techniques for providing strong consistency. Extensive experiments using both a micro-benchmark and the TPC-W benchmark demonstrate that our proposals are viable and achieve considerable scalability while maintaining strong consistency.

Detecting Inconsistencies in Distributed Data

Wenfei Fan, Floris Geerts, Shuai Ma, Heiko Müller; University of Edinburgh, UK

One of the central problems for data quality is inconsistency detection. Given a database D and a set Σ of dependencies as data quality rules, we want to identify tuples in D that violate some rules in Σ . When D is a centralized database, there have been effective SQL-based techniques for finding violations. It is, however, far more challenging when data in D is distributed, in which inconsistency detection often necessarily requires shipping data from one site to another.

This paper develops techniques for detecting violations of conditional functional dependencies (CFDs) in relations that are fragmented and distributed across different sites. (1) We formulate the detection problem in various distributed settings as optimization problems, measured by either network traffic or response time. (2) We show that it is beyond reach in practice to find optimal detection methods: the detection problem is NP-complete when the data is partitioned either horizontally or vertically, and when we aim to minimize either data shipment or response time. (3) For data that is horizontally partitioned, we provide several algorithms to find violations of a set of CFDs, leveraging the structure of CFDs to reduce data shipment or increase parallelism. (4) We verify experimentally that our algorithms are scalable on large relations and complex CFDs. (5) For data that is vertically partitioned, we provide a characterization for CFDs to be checked locally without requiring data shipment, in terms of dependency preservation. We show that it is intractable to minimally refine a partition and make it dependency preserving.

Research Session 3: Stream Mining

Regency D, 10:30 - 12:00, Tuesday Chair: Felix Naumann

Optimal Load Shedding with Aggregates and Mining Queries

Barzan Mozafari, Carlo Zaniolo; University of California at Los Angeles, USA

To cope with bursty arrivals of high-volume data, a DSMS has to shed load while minimizing the degradation of Quality of Service (QoS). In this paper, we show that this problem can be formalized as a classical optimization task from operations research, in ways that accommodate different requirements for multiple users, different query sensitivities to load shedding, and different penalty functions. Standard non-linear programming algorithms are adequate for non-critical situations, but for severe overloads, we propose a more efficient algorithm that runs in linear time, without compromising optimality. Our approach is applicable to a large class of queries including traditional SQL aggregates, statistical aggregates (e.g., quantiles), and data mining functions, such as k-means, naive Bayesian classifiers, decision trees, and frequent pattern discovery (where we can even specify a different error bound for each pattern). In fact, we show that these aggregates, for which the proposed methods apply with full generality.

Finally, we propose a novel architecture for supporting load shedding in an extensible system, where users can write arbitrary User Defined Aggregates (UDA), and thus confirm our analytical findings with several experiments executed on an actual DSMS.

Scheduling for Fast Response Multi-Pattern Matching Over Streaming Events

*Ying Yan*¹, *Jin Zhang*¹, *Ming-Chien Shan*²; ¹*SAP, China*; ²*SAP, USA*

Real-time pattern matching over event streams has gained much more attention recently due to the analytical capability demanded in many operation-critical applications such as credit card fraud detection, algorithmic stock trading and RFID tracking. One of the common but important requirements in the above-mentioned applications is fast response. Usually, there are a large number of pattern queries subscribed in the system, running continuously and concurrently. However, not much research has been done on the scheduling algorithms and management to improve the overall response time of these queries. To address this challenge, we focus on the study of how to improve the average response time of multiple pattern queries. We first propose two static scheduling algorithms: Event-based

(EBS) and Run-based (RBS) Scheduling and discuss what would be a better choice under different system configurations. We then come up with a hybrid method called Fast Response Time Scheduling (FRTS) to dynamically manage the scheduling in order to further reduce the average response time. The experimental results of these scheduling algorithms have shown that the FRTS method can improve 5 times average response time comparing with the basic methods in some cases.

Discovery of Cross-Similarity in Data Streams (Short paper)

Machiko Toyoda, Yasushi Sakurai; NTT, Japan

In this paper, we focus on the problem of finding partial similarity between data streams. Our solution relies on dynamic time warping (DTW) as a similarity measure, which computes the distance between sequences whose lengths and/or sampling rates are different. Instead of straightforwardly using DTW that requires a high computation cost, we propose a streaming method that efficiently detects partial similarity between sequences. Our experiments demonstrate that our method detects pairs of optimal subsequences correctly and that it significantly reduces resources in terms of time and space.

Mining Distribution Change in Stock Order Streams (Short paper)

*Xiaoyan Liu*¹, *Xindong Wu*², *Huaiqing Wang*³, *Rui Zhang*¹, *James Bailey*¹, *Kotagiri Ramamohanarao*¹; ¹*University of Melbourne, Australia*; ²*Hefei University of Technology, China*; ³*City University of Hong Kong, China*

Detecting changes in stock prices is a well known problem in finance with important implications for monitoring and business intelligence. Forewarning of changes in stock price, can be made by the early detection of changes in the distributions of stock order numbers. In this paper, we address the change detection problem for streams of stock order numbers and propose a novel incremental detection algorithm. Our algorithm gains high accuracy and low delay by employing a natural Poisson distribution assumption about the nature of stock order streams. We establish that our algorithm is highly scalable and has linear complexity. We also experimentally demonstrate its effectiveness for detecting change points, via experiments using both synthetic and real-world datasets.

Research Session 4: Location Based Services

Regency B, 13:30 – 15:00, Tuesday Chair: Mohamed Mokbel

TrajStore: An Adaptive Storage System for Very Large Trajectory Data Sets

Philippe Cudre-Mauroux, Eugene Wu, Samuel R. Madden; MIT, USA

The rise of GPS and broadband-speed wireless devices has led to tremendous excitement about a range of applications broadly characterized as "location based services". Current database storage systems, however, are inadequate for manipulating the very large and dynamic spatio-temporal data sets required to support such services. Proposals in the literature either present new indices without discussing how to cluster data, potentially resulting in many disk seeks for lookups of densely packed objects, or use static quadtrees or other partitioning structures, which become rapidly suboptimal as the data or queries evolve. As a result of these performance limitations, we built TrajStore, a dynamic storage system optimized for efficiently retrieving all data in a particular spatiotemporal region. TrajStore maintains an optimal index on the data and dynamically co-locates and compresses spatially and temporally adjacent segments on disk. By letting the storage layer evolve with the index, the system adapts to incoming queries and data and is able to answer most queries via a very limited number of I/Os, even when the queries target regions containing hundreds or thousands of different trajectories.

C3: Concurrency Control on Continuous Queries Over Moving Objects

Jing Dai, Chang-Tien Lu; Virginia Tech, USA

ICDE 2010

Moving object management approaches, especially continuous query processing techniques, have attracted significant research effort due to the broad usage of location-aware devices. However, little attention has been given to designing concurrency control protocols for continuous query processing. Existing concurrency control protocols for spatial indices are based on a single indexing tree, while popular continuous query processing approaches require multiple indices. In addition, continuous monitoring combined with frequent location updates challenges the development of serializable isolation for concurrent index operations. This paper proposes an efficient concurrent continuous query processing approach C3, which fuses scalable continuous query processing methods with lazy update techniques on R-trees. The proposed concurrency control protocol, equipped with intra- and interindex protection, assures serializable isolation, consistency, and deadlock-freedom. The correctness of the proposed protocol is theoretically proven, and the experiment results demonstrated its scalability and efficiency.

Policy-Aware Sender Anonymity in Location Based Services

*Alin Deutsch*¹, *Richard Hull*², *Avinash Vyas*³, *Kevin Keliang Zhao*¹; ¹University of California at San Diego, USA; ²IBM, USA; ³Bell Labs Research, USA

Sender anonymity in location-based services (LBS) attempts to hide the identity of a mobile device user who sends requests to the LBS provider for services in her proximity (e.g. "find the nearest gas station" etc.). The goal is to keep the requester's interests private even from attackers who (via hacking or subpoenas) gain access to the request and to the locations of the mobile user and other nearby users at the time of the request. In an LBS context, the best-studied privacy guarantee is known as *sender k-anonymity*. We show that state-of-the art solutions for sender k-anonymity defend only against naive attackers who have no knowledge of the anonymization policy that is in use. We strengthen the privacy guarantee to defend against more realistic "policy-aware" attackers. We describe a polynomial algorithm to obtain an optimum anonymization policy. Our implementation and experiments show that the policy-aware sender k-anonymity has potential for practical impact, being efficiently enforceable, with limited reduction in utility when compared to policy-unaware guarantees.

Research Session 5: Probabilistic Databases

Regency C, 13:30 – 15:00, Tuesday Chair: Reynold Chang

Approximate Confidence Computation in Probabilistic Databases

Dan Olteanu¹, Jiewen Huang¹, Christoph Koch²; ¹University of Oxford, UK; ²Cornell University, USA

This paper introduces a deterministic approximation algorithm with error guarantees for computing the probability of propositional formulas over discrete random variables. The algorithm is based on an incremental compilation of formulas into decision diagrams using three types of decompositions: Shannon expansion, independence partitioning, and product factorization. With each decomposition step, lower and upper bounds on the probability of the partially compiled formula can be quickly computed and checked against the allowed error.

This algorithm can be effectively used to compute approximate confidence values of answer tuples to positive relational algebra queries on general probabilistic databases (c-tables with discrete probability distributions). We further tune our algorithm so as to capture all known tractable conjunctive queries without self-joins on tuple-independent probabilistic databases: In this case, the algorithm requires time polynomial in the input size even for *exact* computation.

We implemented the algorithm as an extension of the SPROUT query engine. An extensive experimental effort shows that it consistently outperforms state-of-art approximation techniques by several orders of magnitude.

PIP: A Database System for Great and Small Expectations

Oliver Kennedy, Christoph Koch; Cornell University, USA

Estimation via sampling out of highly selective join queries is well known to be problematic, most notably in online aggregation. Without goal-directed sampling strategies, samples falling outside of the selection constraints lower estimation efficiency at best, and cause inaccurate estimates at worst. This problem appears in general probabilistic database systems, where query processing is tightly coupled with sampling. By committing to a set of samples before evaluating the query, the engine wastes effort on samples that will be discarded, query processing that may need to be repeated, or unnecessarily large numbers of samples.

We describe PIP, a general probabilistic database system that uses symbolic representations of probabilistic data to defer computation of expectations, moments, and other statistical measures until the expression to be measured is fully known. This approach is sufficiently general to admit both continuous and discrete distributions. Moreover, deferring sampling enables a broad range of goal-oriented sampling-based (as well as exact) integration techniques for computing expectations, allows the selection of the integration strategy most appropriate to the expression being measured, and can reduce the amount of sampling work required.

We demonstrate the effectiveness of this approach by showing that even straightforward algorithms can make use of the added information. These algorithms have a profoundly positive impact on the efficiency and accuracy of expectation computations, particularly in the case of highly selective join queries.

Generator-Recognizer Networks: A Unified Approach to Probabilistic Databases (Short paper)

Ruiwen Chen, Yongyi Mao, Iluju Kiringa; University of Ottawa, Canada

Under the tuple-level uncertainty paradigm, we introduce a novel graphical model, Generator-Recognizer Network (GRN), as a model for probabilistic databases. The GRN modeling framework extends existing graphical models of probabilistic databases and is capable of representing a much wider range of dependence structures.

Probabilistic Declarative Information Extraction (Short paper)

*Daisy Zhe Wang*¹, *Eirinaios Michelakis*¹, *Michael J. Franklin*¹, *Minos Garofalakis*², *Joseph M. Hellerstein*¹; ¹University of California at Berkeley, USA; ²Technical University of Crete, Greece

Unstructured text represents a large fraction of the world's data. It often contains snippets of structured information (e.g., people's names and zip codes). Information Extraction (IE) techniques identify such structured information in text. In recent years, database research has pursued IE on two fronts: declarative languages and systems for managing IE tasks, and probabilistic databases for querying the output of IE. In this paper, we make the first step to merge these two directions, without loss of statistical robustness, by implementing a state-of-the-art statistical IE model — Conditional Random Fields (CRF) — in the setting of a Probabilistic Database that treats statistical models as first-class data objects. We show that the Viterbi algorithm for CRF inference can be specified declaratively in recursive SQL. We also show the performance benefits relative to a standalone open-source Viterbi implementation. This work opens up the optimization opportunities for queries involving both inference and relational operators over IE models.

Research Session 6: Spatial Indexing

Regency D, 13:30 – 15:00, Tuesday Chair: Cyrus Shahabi

PARINET: A Tunable Access Method for In-Network Trajectories

*Iulian Sandu Popa*¹, *Karine Zeitouni*¹, *Vincent Oria*², *Dominique Barth*¹, *Sandrine Vial*¹; ¹*PRiSM, France;* ²*New Jersey Institute of Technology, USA*

In this paper we propose PARINET, a new access method to efficiently retrieve the trajectories of objects moving in networks. The structure of PARINET is based on a combination of graph partitioning and a set of composite B⁺-tree local indexes. PARINET is designed for historical data and relies on the distribution of the data over the network as for historical data, the data distribution is known in advance. Because the network can be modeled using graphs, the partitioning of the trajectory data is based on graph partitioning theory and can be tuned for a given query load. The data in each partition is indexed on the time component using B⁺-trees. We study different types of queries, and provide an optimal configuration for several scenarios. PARINET can easily be integrated into any RDBMS, which is an essential asset particularly for industrial or commercial applications. The experimental evaluation under an off-the-shelf DBMS shows that PARINET is robust. It also significantly outperforms both MON-tree and another R-tree based access method which are the reference indexing techniques for in-network trajectory databases.

Multi-Guarded Safe Zone: An Effective Technique to Monitor Moving Circular Range Queries

*Muhammad Aamir Cheema*¹, *Ljiljana Brankovic*², *Xuemin Lin*¹, *Wenjie Zhang*¹, *Wei Wang*¹; ¹*University of New South Wales, Australia;* ²*University of Newcastle, Australia*

Given a positive value r, a circular range query returns the objects that lie within the distance r of the query location. In this paper, we study the circular range queries that continuously change their locations. We present an efficient and effective technique to monitor such moving range queries by utilising the concept of a *safe zone*. The safe zone of a query is the area with a property that while the query remains inside it, the results of the query remain unchanged. Hence, the query does not need to be re-evaluated unless it leaves the safe zone. The shape of the safe zone is defined by the so-called *quard objects*. The cost of checking whether a query lies in the safe zone takes k distance computations, where k is the number of the guard objects. Our contributions are as follows. 1) We propose a technique based on powerful pruning rules and a unique access order which efficiently computes the safe zone and minimizes the I/O cost. 2) To show the effectiveness of the safe zone, we theoretically evaluate the probability that a query leaves the safe zone within one time unit and the expected distance a query moves before it leaves the safe zone. Additionally, for the queries that have diameter of the safe zone less than its expected value multiplied by a constant, we also give an upper bound on the expected number of guard objects. This upper bound turns out to be a constant, that is, it does not depend either on the radius r of the query or the density of the objects. The theoretical analysis is verified by extensive experiments. 3) Our thorough experimental study demonstrates that our proposed approach is close to optimal and is an order of magnitude faster than a naïve algorithm.

Geotagging with Local Lexicons to Build Indexes for Textually-Specified Spatial Data

Michael D. Lieberman, Hanan Samet, Jagan Sankaranarayanan; University of Maryland at College Park, USA

The successful execution of location-based and feature-based queries on spatial databases requires the construction of spatial indexes on the spatial attributes. This is not simple when the data is unstructured as is the case when the data is a collection of documents such as news articles, which is the domain of discourse, where the spatial attribute consists of text that can be (but is not required to be) interpreted as the names of locations. In

other words, spatial data is specified using text (known as a *toponym*) instead of geometry, which means that there is some ambiguity involved. The process of identifying and disambiguating references to geographic locations is known as *geotagging* and involves using a combination of internal document structure and external knowledge, including a document-independent model of the audience's vocabulary of geographic locations, termed its *spatial lexicon*. In contrast to previous work, a new spatial lexicon model is presented that distinguishes between a *global lexicon* of locations known to all audiences, and an audience-specific *local lexicon*. Generic methods for inferring audiences' local lexicons are described. Evaluations of this inference method and the overall geotagging procedure indicate that establishing local lexicons cannot be overlooked, especially given the increasing prevalence of highly local data sources on the Internet, and will enable the construction of more accurate spatial indexes.

Research Session 7: Privacy Techniques

Regency B, 15:30 – 17:00, Tuesday Chair: Wei-Shinn Ku

On Optimal Anonymization for l^+ -Diversity

Junqiang Liu, Ke Wang; Simon Fraser University, Canada

Publishing person specific data while protecting privacy is an important problem. Existing algorithms that enforce the privacy principle called *l*-diversity are heuristic based due to the NP-hardness. Several questions remain open: can we get a significant gain in the data utility from an optimal solution compared to heuristic ones; can we improve the utility by setting a distinct privacy threshold per sensitive value; is it practical to find an optimal solution efficiently for real world datasets. This paper addresses these questions. Specifically, we present a pruning based algorithm for finding an *optimal solution* to an extended form of the *l*-diversity problem. The novelty lies in several strong techniques: a novel structure for enumerating all solutions, methods for estimating cost lower bounds, strategies for dynamically arranging the enumeration order and updating lower bounds. This approach can be instantiated with any reasonable cost metric. Experiments on real world datasets show that our algorithm is efficient and improves the data utility.

Differential Privacy via Wavelet Transforms

*Xiaokui Xiao*¹, *Guozhang Wang*², *Johannes Gehrke*²; ¹*Nanyang Technological University, Singapore;* ²*Cornell University, USA*

Privacy preserving data publishing has attracted considerable research interest in recent years. Among the existing solutions, ε -*differential privacy* provides one of the strongest privacy guarantees. Existing data publishing methods that achieve ε -differential privacy, however, offer little data utility. In particular, if the output dataset is used to answer count queries, the noise in the query answers can be proportional to the number of tuples in the data, which renders the results useless.

In this paper, we develop a data publishing technique that ensures ε -differential privacy while providing accurate answers for *range-count queries*, i.e., count queries where the predicate on each attribute is a range. The core of our solution is a framework that applies *wavelet transforms* on the data before adding noise to it. We present instantiations of the proposed framework for both ordinal and nominal data, and we provide a theoretical analysis on their privacy and utility guarantees. In an extensive experimental study on both real and synthetic data, we show the effectiveness and efficiency of our solution.

Efficient Verification of Shortest Path Search via Authenticated Hints

*Man Lung Yiu*¹, *Yimin Lin*², *Kyriakos Mouratidis*²; ¹*Hong Kong Polytechnic University, China*; ²*Singapore Management University, Singapore*

Shortest path search in transportation networks is unarguably one of the most important online search services nowadays (e.g., Google Maps, MapQuest, etc), with applications spanning logistics, spatial optimization, or everyday driving decisions. Often times, the owner

of the road network data (e.g., a transport authority) provides its database to third-party query services, which are responsible for answering shortest path queries posed by their clients. The issue arising here is that a query service might be returning sub-optimal paths either purposely (in order to serve its own purposes like computational savings or commercial reasons) or because it has been compromised by Internet attackers who falsify the results. Therefore, for the above applications to succeed, it is essential that each reported path is accompanied by a proof, which allows clients to verify the path's correctness.

This is the first study on shortest path verification in outsourced network databases. We propose the concept of *authenticated hints*, which is used to reduce the size of the proofs. We develop several authentication techniques and quantify their tradeoffs with respect to offline construction cost and proof size. Experiments on real road networks demonstrate that our solutions are indeed efficient and lead to compact query proofs.

Research Session 8: Skyline Queries

Regency C, 15:30 – 17:00, Tuesday Chair: Xuemin Lin

Evaluating Skylines in the Presence of Equijoins

*Wen Jin*¹, *Michael D. Morse*¹, *Jignesh M. Patel*², *Martin Ester*³, *Zengjian Hu*³; ¹*University of Michigan, USA*; ²*University of Wisconsin-Madison, USA*; ³*Simon Fraser University, Canada*

When a database system is extended with the skyline operator, it is important to determine the most efficient way to execute a skyline query across tables with join operations. This paper describes a framework for evaluating skylines in the presence of equijoins, including: (1) the development of algorithms to answer such queries over large input tables in a non-blocking, pipeline fashion, which significantly speeds up the entire query evaluation time. These algorithms are built on top of the traditional relational *Nested-Loop* and the *Sort-Merge* join algorithms, which allows easy implementation of these methods in existing relational systems; (2) a novel method for estimating the skyline selectivity of the joined table; (3) evaluation techniques; and (4) a systematic experimental evaluation to validate our skyline evaluation framework.

Route Skyline Queries: A Multi-Preference Path Planning Approach

Hans-Peter Kriegel, Matthias Renz, Matthias Schubert; LMU München, Germany

In recent years, the research community introduced various methods for processing skyline queries in multidimensional databases. The skyline operator retrieves all objects being optimal w.r.t. an arbitrary linear weighting of the underlying criteria. The most prominent example query is to find a reasonable set of hotels which are cheap but close to the beach. In this paper, we propose a new approach for computing skylines on routes (paths) in a road network considering multiple preferences like *distance*, *driving time*, *the number of traffic* lights, gas consumption, etc. Since the consideration of different preferences usually involves different routes, a skyline-fashioned answer with relevant route candidates is highly useful. In our work, we employ graph embedding techniques to enable a best-first based graph exploration considering route preferences based on arbitrary road attributes. The core of our skyline query processor is a route iterator which iteratively computes the top routes according to (at least one) preference in an efficient way avoiding that route computations need to be issued from scratch in each iteration. Furthermore, we propose pruning techniques in order to reduce the search space. Our pruning strategies aim at pruning as many route candidates as possible during the graph exploration. Therefore, we are able to prune candidates which are only partially explored. Finally, we show that our approach is able to reduce the search space significantly and that the skyline can be computed in efficient time in our experimental evaluation.

Probabilistic Contextual Skylines

Dimitris Sacharidis, Anastasios Arvanitis, Timos Sellis; Athena RC, Greece

The skyline query returns the most interesting tuples according to a set of explicitly defined preferences among attribute values. This work relaxes this requirement, and allows users to pose meaningful skyline queries without stating their choices. To compensate for missing knowledge, we first determine a set of uncertain preferences based on user profiles, i.e., information collected for previous contexts. Then, we define a probabilistic contextual skyline query (p-CSQ) that returns the tuples which are interesting with high probability. We emphasize that, unlike past work, uncertainty lies within the query and not the data, i.e., it is in the relationships among tuples rather than in their attribute values. Furthermore, due to the nature of this uncertainty, popular skyline methods, which rely on a particular tuple visit order, do not apply for p-CSQs. Our experimental evaluation concludes that the proposed techniques are significantly more efficient compared to a standard block nested loops approach.

Research Session 9: Information Integration

Regency D, 15:30 - 17:00, Tuesday Chair: Mourad Ouzzani

Schema Covering: A Step Towards Enabling Reuse in Information Integration

*Barna Saha*¹, *Ioana Stanoi*², *Kenneth L. Clarkson*²; ¹*University of Maryland at College Park, USA*; ²*IBM, USA*

We introduce *schema covering*, the problem of identifying easily understandable common objects for describing large and complex schemas. Defining transformations between schemas is a key objective in information integration. However, this process often becomes cumbersome when the schemas are large and structurally complex. If such complex schemas can be broken into smaller and simpler objects, then simple transformations defined over these smaller objects can be reused to define suitable transformations among the complex schemas. Schema covering performs this vital task by identifying a collection of common concepts from a repository and creating a *cover* of the complex schema by these concepts. In this paper, we formulate the problem of schema covering, show that it is NP-Complete, and give efficient approximation algorithms for it. A performance evaluation with real business schemas confirms the effectiveness of our approach.

Managing Uncertainty of XML Schema Matching

Reynold Cheng, Jian Gong, David W. Cheung; University of Hong Kong, China

Despite of advances in machine learning technologies, a schema matching result between two database schemas (e.g., those derived from COMA++) is likely to be imprecise. In particular, numerous instances of "possible mappings" between the schemas may be derived from the matching result. In this paper, we study the problem of managing possible mappings between two heterogeneous XML schemas. We observe that for XML schemas, their possible mappings have a high degree of overlap. We hence propose a novel data structure, called the *block tree*, to capture the commonalities among possible mappings. The block tree is useful for representing the possible mappings in a compact manner, and can be generated efficiently. Moreover, it supports the evaluation of *probabilistic twig query* (PTQ), which returns the probability of portions of an XML document that match the query pattern. For users who are interested only in answers with *k*-highest probabilities, we also propose the top-*k* PTQ, and present an efficient solution for it.

The second challenge we have tackled is to efficiently generate possible mappings for a given schema matching. While this problem can be solved by existing algorithms, we show how to improve the performance of the solution by using a divide-and-conquer approach. An extensive evaluation on realistic datasets show that our approaches significantly improve the efficiency of generating, storing, and querying possible mappings.

Propagating Updates Through XML Views Using Lineage Tracing

Leonidas Fegaras; University of Texas at Arlington, USA

We address the problem of updating XML views over relational data by translating view updates expressed in the XQuery update facility to embedded SQL updates. Although our XML views may be defined using the full extent of the XQuery syntax, they can only connect relational tables through restricted one-to-many relationships that do not cause view side effects for a wide range of XQuery updates. Our approach is to use lineage tracing to propagate the necessary information about the origins of updatable data pieces through the guery and the view code, to be used when these pieces are to be updated. Our system performs a compile-time analysis, based on polymorphic type inference and type usage, to detect the exclusive data sources, which are the table columns from the database that can be updated without causing side-effects to the view. The rest of the updates are associated with an update context in the form of a chain of tuples, which reflects the navigation path that was used to reach the update destination. At commit time, our system collectively considers all the compatible chains of all updates in the transaction and tries to relink them to new chains from the existing database whose tuples contain the updated data, so that the updates are reflected correctly without causing side effects to the other components of the view.

Research Session 10: Query Interfaces

Beacon A, 15:30 – 17:00, Tuesday Chair: Panos Ipeirotis

USHER: Improving Data Quality with Dynamic Forms

*Kuang Chen*¹, *Harr Chen*², *Neil Conway*¹, *Joseph M. Hellerstein*¹, *Tapan S. Parikh*¹; ¹*University of California at Berkeley, USA*; ²*MIT, USA*

Data quality is a critical problem in modern databases. Data entry forms present the first and arguably best opportunity for detecting and mitigating errors, but there has been little research into automatic methods for improving data quality at entry time. In this paper, we propose USHER, an end-to-end system for form design, entry, and data quality assurance. Using previous form submissions, USHER learns a probabilistic model over the questions of the form. USHER then applies this model at every step of the data entry process to improve data quality. Before entry, it induces a form layout that captures the most important data values of a form instance as quickly as possible. During entry, it dynamically adapts the form to the values being entered, and enables real-time feedback to guide the data enterer toward their intended values. After entry, it re-asks questions that it deems likely to have been entered incorrectly. We evaluate all three components of USHER using two real-world data sets. Our results demonstrate that each component has the potential to improve data quality considerably, at a reduced cost when compared to current practice.

Explaining Structured Queries in Natural Language

*Georgia Koutrika*¹, *Alkis Simitsis*², *Yannis E. Ioannidis*³; ¹*Stanford University, USA*; ²*HP, USA*; ³*University of Athens, Greece*

Many applications offer a form-based environment for naïve users for accessing databases without being familiar with the database schema or a structured query language. User interactions are translated to structured queries and executed. However, as a user is unlikely to know the underlying semantic connections among the fields presented in a form, it is often useful to provide her with a textual explanation of the query. In this paper, we take a graph-based approach to the query translation problem. We represent various forms of structured queries as directed graphs and we annotate the graph edges with template labels using an extensible template mechanism. We present different graph traversal strategies for efficiently exploring these graphs and composing textual query descriptions. Finally, we present experimental results for the efficiency and effectiveness of the proposed methods.

ScoreFinder: A Method for Collaborative Quality Inference on User-Generated Content (Short paper)

Yang Liao, Aaron Harwood, Kotagiri Ramamohanarao; University of Melbourne, Australia

User-generated content is quickly becoming the greatest source of information on the World Wide Web. Shared content items are initially considered *unconfirmed* in the sense that their *credibility* has not yet been established. Conventional, centralized confirmation of credibility is infeasible at the Internet scale and so making use of the *annotators* themselves to evaluate each item is essential. However, users usually differ in opinions to the same item, and the existence of bias, variance and maliciousness makes the problem of aggregating opinions more difficult. Addressing this problem, we propose the use of an *Author-Annotator model* with an iterative algorithm, called *ScoreFinder*, for inferring credibility by ranking shared items. In order to reduce the influence from a variety of error sources, we identify reliable users on each topic, and adaptively aggregate scores from them. Moreover, we transform the users' input to remove errors/anomalies, by identifying patterns of misbehaviour learned from a real data sets, and a significant improvement was achieved in the experiment.

IQ^{*P*}: **Incremental Query Construction, a Probabilistic Approach** (Short paper)

*Elena Demidova*¹, *Xuan Zhou*², *Wolfgang Nejdl*¹; ¹L3S Research Center, *Germany*; ²CSIRO, Australia

This paper presents IQ^P — a novel approach to bridge the gap between usability of keyword search and expressiveness of database queries. IQ^P enables a user to start with an arbitrary keyword query and incrementally refine it into a structured query through an interactive interface. The enabling techniques of IQ^P include: (1) a conceptual framework for incremental query construction; (2) a probabilistic model to assess the possible informational needs represented by a keyword query; (3) an algorithm to perform an optimal query construction.

Research Session 11: Top-K Queries

Regency B, 10:30 – 12:00, Wednesday Chair: Ralf Schenkel

TASM: Top-k Approximate Subtree Matching

Nikolaus Augsten¹, Denilson Barbosa², Michael Böhlen¹, Themis Palpanas³; ¹Free University of Bozen-Bolzano, Italy; ²University of Alberta, Canada; ³University of Trento, Italy

We consider the *Top-k Approximate Subtree Matching* (TASM) problem: finding the *k* best matches of a small query tree, e.g., a DBLP article with 15 nodes, in a large document tree, e.g., DBLP with 26M nodes, using the canonical tree edit distance as a similarity measure between subtrees. Evaluating the tree edit distance for large XML trees is difficult: the best known algorithms have cubic runtime and quadratic space complexity, and, thus, do not scale. Our solution is TASM-postorder, a memory-efficient and scalable TASM algorithm. We prove an upper-bound for the maximum subtree size for which the tree edit distance needs to be evaluated. The upper bound depends on the query and is independent of the document size and structure. A core problem is to efficiently prune subtrees that are above this size threshold. We develop an algorithm based on the prefix ring buffer that allows us to prune all subtrees above the threshold in a single postorder scan of the document. The size of the prefix ring buffer is linear in the threshold. As a result, the space complexity of TASM-postorder depends only on *k* and the query size, and the runtime of TASM-postorder is linear in the size of the document. Cure space complexity of TASM-postorder depends only on *k* and the query size, and the runtime of TASM-postorder is linear in the size of the documents confirms our analytic results.

Reverse Top-k Queries

*Akrivi Vlachou*¹, *Christos Doulkeridis*¹, *Yannis Kotidis*², *Kjetil Nørvåg*¹; ¹*NTNU*, *Norway*; ²*AUEB*, *Greece*

Rank-aware query processing has become essential for many applications that return to the user only the top-*k* objects based on the individual user's preferences. Top-*k* queries have been mainly studied from the perspective of the user, focusing primarily on efficient query processing. In this work, for the first time, we study top-k queries from the perspective of the product manufacturer. Given a potential product, which are the user preferences for which this product is in the top-k query result set? We identify a novel query type, namely reverse top-k query, that is essential for manufacturers to assess the potential market and impact of their products based on the competition. We formally define reverse top-k queries and introduce two versions of the query, namely monochromatic and bichromatic. We first provide a geometric interpretation of the monochromatic reverse top-kquery in the solution space that helps to understand the reverse top-k query conceptually. Then, we study in more details the case of bichromatic reverse top-k query, which is more interesting for practical applications. Such a query, if computed in a straightforward manner, requires evaluating a top-k query for each user preference in the database, which is prohibitively expensive even for moderate datasets. In this paper, we present an efficient threshold-based algorithm that eliminates candidate user preferences, without processing the respective top-k queries. Furthermore, we introduce an indexing structure based on materialized reverse top-k views in order to speed up the computation of reverse top-k queries. Materialized reverse top-k views trade preprocessing cost for query speed up in a controllable manner. Our experimental evaluation demonstrates the efficiency of our techniques, which reduce the required number of top-k computations by 1 to 3 orders of magnitude.

Top-K Aggregation Queries Over Large Networks (Short paper)

*Xifeng Yan*¹, *Bin He*², *Feida Zhu*³, *Jiawei Han*⁴; ¹*University of California at Santa Barbara, USA*; ²*IBM, USA*; ³*Singapore Management University, Singapore*; ⁴*University of Illinois at Urbana-Champaign, USA*

Searching and mining large graphs today is critical to a variety of application domains, ranging from personalized recommendation in social networks, to searches for functional associations in biological pathways. In these domains, there is a need to perform aggregation operations on large-scale networks. Unfortunately the existing implementation of aggregation operations on relational databases does not guarantee superior performance in network space, especially when it involves edge traversals and joins of gigantic tables.

In this paper, we investigate the neighborhood aggregation queries: Find nodes that have top-k highest aggregate values over their h-hop neighbors. While these basic queries are common in a wide range of search and recommendation tasks, surprisingly they have not been studied systematically. We developed a Local Neighborhood Aggregation framework, called LONA, to answer them efficiently. LONA exploits two properties unique in network space: First, the aggregate value for the neighboring nodes should be similar in most cases; Second, given the distribution of attribute values, it is possible to estimate the upper-bound value of aggregates. These two properties inspire the development of novel pruning techniques, forward pruning using differential index and backward pruning using partial distribution. Empirical results show that LONA could outperform the baseline algorithm up to 10 times in real-life large networks.

TopCells: Keyword-Based Search of Top-k Aggregated Documents in Text Cube (Short paper)

Bolin Ding, Bo Zhao, Cindy Xide Lin, Jiawei Han, Chengxiang Zhai; University of Illinois at Urbana-Champaign, USA

Previous studies on supporting keyword queries in RDBMSs provide users with a ranked list of relevant linked structures (*e.g.* joined tuples) or individual tuples. In this paper, we aim to support keyword search in a data cube with text-rich dimension(s) (so-called *text cube*). Each document is associated with structural dimensions. A cell in the text cube aggregates

a set of documents with matching dimension values on a subset of dimensions. Given a keyword query, our goal is to find the top-k most relevant cells in the text cube. We propose a relevance scoring model and efficient ranking algorithms. Experiments are conducted to verify their efficiency.

Research Session 12: Workflow and Workload Management

Regency C, 10:30 – 12:00, Wednesday Chair: Holger Schwarz

Optimizing ETL Workflows for Fault-Tolerance

Alkis Simitsis, Kevin Wilkinson, Umeshwar Dayal, Malu Castellanos; HP, USA

Extract-Transform-Load (ETL) processes play an important role in data warehousing. Typically, design work on ETL has focused on performance as the sole metric to make sure that the ETL process finishes within an allocated time window. However, other quality metrics are also important and need to be considered during ETL design. In this paper, we address ETL design for performance plus fault-tolerance and freshness. There are many reasons why an ETL process can fail and a good design needs to guarantee that it can be recovered within the ETL time window. How to make ETL robust to failures is not trivial. There are different strategies that can be used and they each have different costs and benefits. In addition, other metrics can affect the choice of a strategy; e.g., higher freshness reduces the time window for recovery. The design space is too large for informal, ad-hoc approaches. In this paper, we describe our QoX optimizer that considers multiple design strategies and finds an ETL design that satisfies multiple objectives. In particular, we define the optimizer search space, cost functions, and search algorithms. Also, we illustrate its use through several experiments and we show that it produces designs that are very near optimal.

Q-Cop: Avoiding Bad Query Mixes to Minimize Client Timeouts Under Heavy Loads

Sean Tozer, Tim Brecht, Ashraf Aboulnaga; University of Waterloo, Canada

In three-tiered web applications, some form of admission control is required to ensure that throughput and response times are not significantly harmed during periods of heavy load. We propose Q-Cop, a prototype system for improving admission control decisions that considers a combination of the load on the system, the number of simultaneous queries being executed, the actual mix of queries being executed, and the expected time a user may wait for a reply before they or their browser give up (i.e., time out). Using TPC-W queries, we show that the response times of different types of queries can vary significantly depending not just on the number of queries being processed but on the mix of other queries that accounts for the mix of queries being executed and integrate this model into a three-tiered system to make admission control decisions. Our results show that this approach makes more informed decisions about which queries to reject and as a result significantly reduces the number of requests that time out. Across the range of workloads examined an average of 47% fewer requests are unsuccessful than the next best approach.

Admission Control Mechanisms for Continuous Queries in the Cloud (Short paper)

Lory Al Moakar¹, Panos K. Chrysanthis¹, Christine Chung², Shenoda Guirguis¹, Alexandros Labrinidis¹, Panayiotis Neophytou¹, Kirk Pruhs¹; ¹University of Pittsburgh, USA; ²Connecticut College, USA

Amazon, Google, and IBM now sell cloud computing services.We consider the setting of a for-profit business selling data stream monitoring/management services and we investigate auction-based mechanisms for admission control of continuous queries. When submitting a query, each user also submits a bid of how much she is willing to pay for that query to run. The admission control auction mechanism then determines which queries to admit,

and how much to charge each user in a way that maximizes system revenue while being strategyproof and sybil immune, incentivizing users to use the system honestly. Specifically, we require that each user maximizes her payoff by bidding her true value of having her query run. We design several payment mechanisms and experimentally evaluate them. We describe the provable game theoretic characteristics of each mechanism alongside its performance with respect to maximizing profit and total user payoff.

Interaction-Aware Prediction of Business Intelligence Workload Completion Times (Short paper)

*Mumtaz Ahmad*¹, *Songyun Duan*², *Ashraf Aboulnaga*¹, *Shivnath Babu*²; ¹*University of Waterloo, Canada;* ²*Duke University, USA*

While planning the execution of report-generation workloads, database administrators often need to know how long different query workloads will take to run. Database systems run mixes of multiple queries of different types concurrently. Hence, estimating the completion time of a query workload requires reasoning about query mixes and inter-query interactions in the mixes; rather than considering queries or query types in isolation. This paper presents a novel approach for estimating workload completion time based on experimentdriven modeling and simulation of the impact of inter-query interactions. A preliminary evaluation of this approach with TPC-H queries on IBM DB2 shows how our approach can consistently predict workload completion times with good accuracy.

Research Session 13: Indexing and Hashing

Regency D, 10:30 – 12:00, Wednesday Chair: Paul Larson

Fast In-Memory XPath Search Using Compressed Indexes

Diego Arroyuelo¹, Francisco Claude², Sebastian Maneth³, Veli Mäkinen⁴, Gonzalo Navarro⁵, Kim Nguyễn³, Jouni Sirén⁴, Niko Välimäki⁴; ¹Yahoo!, Chile; ²University of Waterloo, Canada; ³NICTA, Australia; ⁴University of Helsinki, Finland; ⁵University of Chile, Chile

A large fraction of an XML document typically consists of text data. The XPath query language allows text search via the equal, contains, and starts-with predicates. Such predicates can be efficiently implemented using a compressed self-index of the document's text nodes. Most queries, however, contain some parts querying the text of the document, plus some parts querying the tree structure. It is therefore a challenge to choose an appropriate evaluation order for a given query, which optimally leverages the execution speeds of the text and tree indexes. Here the SXSI system is introduced. It stores the tree structure of an XML document using a bit array of opening and closing brackets plus a sequence of labels, and stores the text nodes of the document using a global compressed self-index. On top of these indexes sits an XPath query engine that is based on tree automata. The engine uses fast counting queries of the text index in order to dynamically determine whether to evaluate top-down or bottom-up with respect to the tree structure. The resulting system has several advantages over existing systems: (1) on pure tree queries (without text search) such as the XPathMark queries, the SXSI system performs on par or better than the fastest known systems MonetDB and Qizx, (2) on queries that use text search, SXSI outperforms the existing systems by 1-3 orders of magnitude (depending on the size of the result set), and (3) with respect to memory consumption, SXSI outperforms all other systems for countingonly queries.

Hashing Tree-Structured Data: Methods and Applications

Shirish Tatikonda, Srinivasan Parthasarathy; Ohio State University, USA

In this article we propose a new hashing framework for tree-structured data. Our method maps an unordered tree into a multiset of simple wedge-shaped structures refered to as *pivots*. By coupling our pivot multisets with the idea of minwise hashing, we realize a fixed sized signature-sketch of the tree-structured datum yielding an effective mechanism

for hashing such data. We discuss several potential pivot structures and study some of the theoretical properties of such structures, and discuss their implications to tree edit distance and properties related to perfect hashing. We then empirically demonstrate the efficacy and efficiency of the overall approach on a range of real-world datasets and applications.

Estimating the Compression Fraction of an Index Using Sampling (Short paper)

*Stratos Idreos*¹, *Raghav Kaushik*², *Vivek Narasayya*², *Ravishankar Ramamurthy*²; ¹*CWI*, *The Netherlands*; ²*Microsoft*, *USA*

Data compression techniques such as null suppression and dictionary compression are commonly used in today's database systems. In order to effectively leverage compression, it is necessary to have the ability to efficiently and accurately estimate the size of an index if it were to be compressed. Such an analysis is critical if automated physical design tools are to be extended to handle compression. Several database systems today provide estimators for this problem based on random sampling. While this approach is efficient, there is no previous work that analyses its accuracy. In this paper, we analyse the problem of estimating the compressed size of an index from the point of view of worst-case guarantees. We show that the simple estimator implemented by several database systems has several "good" cases even though the estimator itself is agnostic to the internals of the specific compression algorithm.

The Hybrid-Layer Index: A Synergic Approach to Answering Top-*k* Queries in Arbitrary Subspaces (Short paper)

*Jun-Seok Heo*¹, *Junghoo Cho*², *Kyu-Young Whang*¹; ¹*KAIST, Korea*; ²*University of California at Los Angeles, USA*

In this paper, we propose the *Hybrid-Layer Index* (simply, the *HL-index*) that is designed to answer top-*k* queries efficiently when the queries are expressed on any *arbitrary subset* of attributes in the database. Compared to existing approaches, the HL-index significantly reduces the number of tuples accessed during query processing by pruning unnecessary tuples based on two criteria, i.e., it filters out tuples both (1) *globally* based on the combination of *all* attribute values of the tuples like in the layer-based approach (simply, *layer-level filtering*) and (2) based on *individual* attribute values used for ranking the tuples like in the list-based approach (simply, *list-level filtering*). Specifically, the HL-index exploits the synergic effect of integrating the layer-level filtering method and the list-level filtering method. Details and extensive experiments are available in the full paper [7].

Research Session 14: Scientific Data Mining

Regency B, 14:00 – 15:30, Wednesday Chair: Ambuj Singh

The Model-Summary Problem and a Solution for Trees

*Biswanath Panda*¹, *Mirek Riedewald*², *Daniel Fink*³; ¹*Google, USA*; ²*Northeastern University, USA*; ³*Cornell University, USA*

Modern science is collecting massive amounts of data from sensors, instruments, and through computer simulation. It is widely believed that analysis of this data will hold the key for future scientific breakthroughs. Unfortunately, deriving knowledge from large high-dimensional scientific datasets is difficult. One emerging answer is exploratory analysis using data mining; but data mining models that accurately capture natural processes tend to be very complex and are usually not intelligible. Scientists therefore generate model summaries to find the most important patterns learned by the model. We formalize the model-summary problem and introduce it as a novel problem to the database community. Generating model summaries creates serious data management challenges: Scientists usually want to analyze patterns in different "slices" and "dices" of the data space, comparing the effects of various input variables on the output. We propose novel techniques for efficiently generating such summaries for the popular class of tree-based models. Our techniques leverage workload structure on multiple levels. We also propose a scalable implementation of our techniques in MapReduce. For both sequential and parallel implementation, we achieve speedups of one or more orders of magnitude over the naive algorithm, while guaranteeing the exact same results.

Efficient and Accurate Discovery of Patterns in Sequence Datasets

Avrilia Floratou¹, Sandeep Tata², Jignesh M. Patel¹; ¹University of Wisconsin-Madison, USA; ²IBM, USA

Existing sequence mining algorithms mostly focus on mining for subsequences. However, a large class of applications, such as biological DNA and protein motif mining, require efficient mining of "approximate" patterns that are contiguous. The few existing algorithms that can be applied to find such contiguous approximate pattern mining have drawbacks like poor scalability, lack of guarantees in finding the pattern, and difficulty in adapting to other applications. In this paper, we present a new algorithm called FLAME (FLexible and Accurate Motif DEtector). FLAME is a flexible suffix tree based algorithm that can be used to find frequent patterns with a variety of definitions of motif (pattern) models. It is also accurate, as it always find the pattern if it exists. Using both real and synthetic datasets, we demonstrate that FLAME is fast, scalable, and outperforms existing algorithms on a variety of performance metrics. Using FLAME, it is now possible to mine datasets that would have been prohibitively difficult with existing tools.

Mining Mutation Chains in Biological Sequences

*Chang Sheng*¹, *Wynne Hsu*¹, *Mong Li Lee*¹, *Joo Chuan Tong*², *See-Kiong Ng*²; ¹*National University of Singapore, Singapore;* ²*Institute of Infocomm Research, Singapore*

The increasing infectious disease outbreaks has led to a need for new research to better understand the disease's origins, epidemiological features and pathogenicity caused by fast-mutating, fast-spreading viruses. Traditional sequence analysis methods do not take into account the spatio-temporal dynamics of rapidly evolving and spreading viral species. They are also focused on identifying single-point mutations. In this paper, we propose a novel approach that incorporates space-time relationships for studying changes in protein sequences from fast mutating viruses. We aim to detect both single-point mutations as well as k-mutations in the viral sequences. We define the problem of mutation chain pattern mining and design algorithms to discover valid mutation chains. Compact data structures to facilitate the mining process as well as pruning strategies to increase the scalability of the algorithms are devised. Experiments on both synthetic datasets and real world influenza A virus dataset show that our algorithms are scalable and effective in discovering mutations that occur geographically over time.

Research Session 15: Database Performance and Reliability

Regency C, 14:00 – 15:30, Wednesday Chair: Shahin Shayandeh

Exploring Power-Performance Tradeoffs in Database Systems

*Zichen Xu*¹, *Yi-Cheng Tu*¹, *Xiaorui Wang*²; ¹*University of South Florida, USA*; ²*University of Tennessee, USA*

With the total energy consumption of computing systems increasing in a steep rate, much attention has been paid to the design of energy-efficient computing systems and applications. So far, database system design has focused on improving performance of query processing. The objective of this study is to experimentally explore the potential of power conservation in relational database management systems. We hypothesize that, by modifying the query optimizer in a DBMS to take the power cost of query plans into consideration, we will be able to reduce the power usage of database servers and control the tradeoffs between power consumption and system performance. We also identify the sources of such savings by investigating the resource consumption features during query processing

in DBMSs. To that end, we provide an in-depth anatomy and qualitatively analyze the power profile of typical queries in the TPC benchmarks. We perform extensive experiments on a physical testbed based on the PostgreSQL system using workloads generated from the TPC benchmarks. Our hypothesis is supported by such experimental results: power savings in the range of 11%–22% can be achieved by equipping the DBMS with a query optimizer that selects query plans based on both estimated processing time and power requirements.

Workload Driven Index Defragmentation

Vivek Narasayya, Manoj Syamala; Microsoft, USA

Decision support queries that scan large indexes can suffer significant degradation in I/O performance due to index fragmentation. DBAs rely on rules of thumb that use index size and fragmentation information to accomplish the task of deciding which indexes to defragment. However, there are two fundamental limitations that make this task challenging. First, database engines offer little support to help estimate the impact of defragmenting an index on the I/O performance of a query. Second, defragmentation is supported only at the granularity of an entire B+-Tree, which can be too restrictive since defragmentation is an expensive operation. This paper describes techniques for addressing the above limitations. We also study the problem of selecting the appropriate indexes to defragment for a given workload. We have implemented our techniques in Microsoft SQL Server and developed a tool that can provide appropriate index defragmentation recommendations to DBAs. We evaluate the effectiveness of the proposed techniques on several real and synthetic databases.

Impact of Disk Corruption on Open-Source DBMS

Sriram Subramanian, Yupu Zhang, Rajiv Vaidyanathan, Haryadi S. Gunawi, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, Jeffrey F. Naughton; University of Wisconsin-Madison, USA

Despite the best intentions of disk and RAID manufacturers, on-disk data can still become corrupted. In this paper, we examine the effects of corruption on database management systems. Through injecting faults into the MySQL DBMS, we find that in certain cases, corruption can greatly harm the system, leading to untimely crashes, data loss, or even in-correct results. Overall, of 145 injected faults, 110 lead to serious problems. More detailed observations point us to three deficiencies: MySQL does not have the capability to detect some corruptions due to lack of redundant information, does not isolate corrupted data from valid data, and has inconsistent reactions to similar corruption scenarios.

To detect and repair corruption, a DBMS is typically equipped with an offline checker. Unfortunately, the MySQL offline checker is not comprehensive in the checks it performs, misdiagnosing many corruption scenarios and missing others. Sometimes the checker itself crashes; more ominously, its incorrect checking can lead to incorrect repairs. Overall, we find that the checker does not behave correctly in 18 of 145 injected corruptions, and thus can leave the DBMS vulnerable to the problems described above.

Research Session 16: Spatial Databases

Regency D, 14:00 – 15:30, Wednesday Chair: Michael Böhlen

Locating Mapped Resources in Web 2.0

Dongxiang Zhang, Beng Chin Ooi, Anthony K.H. Tung; National University of Singapore, Singapore

Mapping mashups are emerging Web 2.0 applications in which data objects such as blogs, photos and videos from different sources are combined and marked in a map using APIs that are released by online mapping solutions such as Google and Yahoo Maps. These objects are typically associated with a set of tags capturing the embedded semantic and a set of coordinates indicating their geographical locations. Traditional web resource searching strategies are not effective in such an environment due to the lack of the gazetteer context

in the tags. Instead, a better alternative approach is to locate an object by tag matching. However, the number of tags associated with each object is typically small, making it difficult for an object to capture the complete semantics in the query objects.

In this paper, we focus on the fundamental application of locating geographical resources and propose an efficient tag-centric query processing strategy. In particular, we aim to find a set of nearest co-located objects which together match the query tags. Given the fact that there could be large number of data objects and tags, we develop an efficient search algorithm that can scale up in terms of the number of objects and tags. Further, to ensure that the results are relevant, we also propose a geographical context sensitive *geo-tf-idf* ranking mechanism. Our experiments on synthetic data sets demonstrate its scalability while the experiments using the real life data set confirm its practicality.

Preference Queries in Large Multi-Cost Transportation Networks

*Kyriakos Mouratidis*¹, *Yimin Lin*¹, *Man Lung Yiu*²; ¹*Singapore Management University, Singapore*; ²*Hong Kong Polytechnic University, China*

Research on spatial network databases has so far considered that there is a single cost value associated with each road segment of the network. In most real-world situations, however, there may exist multiple cost types involved in transportation decision making. For example, the different costs of a road segment could be its Euclidean length, the driving time, the walking time, possible toll fee, etc. The relative significance of these cost types may vary from user to user. In this paper we consider such *multi-cost transportation networks* (MCN), where each edge (road segment) is associated with multiple cost values. We formulate skyline and top-*k* queries in MCNs and design algorithms for their efficient processing. Our solutions have two important properties in preference-based querying; the skyline methods are *progressive* and the top-*k* ones are *incremental*. The performance of our techniques is evaluated with experiments on a real road network.

Approximate String Search in Spatial Databases

*Bin Yao*¹, *Feifei Li*¹, *Marios Hadjieleftheriou*², *Kun Hou*¹; ¹*Florida State University, USA*; ²*AT&T Labs Research, USA*

This work presents a novel index structure, MHR-tree, for efficiently answering approximate string match queries in large spatial databases. The MHR-tree is based on the R-tree augmented with the min-wise signature and the linear hashing technique. The min-wise signature for an index node u keeps a concise representation of the union of q-grams from strings under the sub-tree of u. We analyze the pruning functionality of such signatures based on set resemblance between the query string and the q-grams from the sub-trees of index nodes. MHR-tree supports a wide range of query predicates efficiently, including range and nearest neighbor queries. We also discuss how to estimate range query selectivity accurately. We present a novel adaptive algorithm for finding balanced partitions using both the spatial and string information stored in the tree. Extensive experiments on large real data sets demonstrate the efficiency and effectiveness of our approach.

Research Session 17: Sensor Networks

Regency B, 16:00 – 17:30, Wednesday Chair: Farnoush Banaei-Kashani

Global Iceberg Detection Over Distributed Data Streams

*Haiquan Zhao*¹, *Ashwin Lall*¹, *Mitsunori Ogihara*², *Jun Xu*¹; ¹*Georgia Institute of Technology, USA*; ²*University of Miami, USA*

In today's Internet applications or sensor networks we often encounter large amounts of data spread over many physically distributed nodes. The sheer volume of the data and bandwidth constraints make it impractical to send all the data to one central node for query processing. Finding distributed icebergs — elements that may have low frequency

at individual nodes but high aggregate frequency — is a problem that arises commonly in practice. In this paper we present a novel algorithm with two notable properties. First, its accuracy guarantee and communication cost are independent of the way in which element counts (for both icebergs and non-icebergs) are split amongst the nodes. Second, it works even when each distributed data set is a stream (i.e., one pass data access only).

Our algorithm builds upon sketches constructed for the estimation of the second frequency moment (F_2) of data streams. The intuition of our idea is that when there are global icebergs in the union of these data streams the F_2 of the union becomes very large. This quantity can be estimated due to the summable nature of F_2 sketches. Our key innovation here is to establish tight theoretical guarantees of our algorithm, under certain reasonable assumptions, using an interesting combination of convex ordering theory and large deviation techniques.

Non-Dyadic Haar Wavelets for Streaming and Sensor Data

Chetan Gupta, Choudur Lakshminarayan, Song Wang, Abhay Mehta; HP, USA

In streaming and sensor data applications, the problems of synopsis construction and outlier detection are important. Due to their low complexity, desirable properties and relative ease of understanding, wavelet based techniques are often used for both synopsis construction and anomaly detection. In streaming data literature, Mallat's algorithm [1] is often used to achieve a Haar wavelet decomposition in O(n) time. However, there is one limitation to this popular technique, in that it leads to a *dyadic decomposition* of data.

We demonstrate that the property of *non-dyadicity* is of considerable use in synopsis construction and anomaly detection. In this regard we present several application results, a synopsis data structure for streaming data that is an order of magnitude superior to the popular Haar based wavelet technique, a method for finding anomalies for sensor data over non-dyadic hierarchies, etc. In our work, we enable non-dyadicity by proposing a Mallat like construction for a wavelet system that admits non-dyadic basis. Our algorithm builds a non-dyadic hierarchical structure, and is more efficient than the state of the art construction. We prove the correctness of our construction by showing that our basis functions demonstrates the properties of a wavelet system.

Ratio Threshold Queries Over Distributed Data Sources (Short paper)

*Rajeev Gupta*¹, *Krithi Ramamritham*², *Mukesh Mohania*¹; ¹*IBM, India*; ²*IIT Bombay, India*

In this paper we consider triggers over distributed data from various sources such as: "Notify when sale of luxury goods constitute more than 20% of the overall sales". In such queries client desires to be notified whenever the ratio of two aggregates, over distributed data, crosses the specified threshold. The challenge lies in being able to execute the queries with the minimal amount of communication necessary for update propagation. We address the challenge by proposing schemes for converting the client threshold condition into conditions on individual distributed data sources such that (1) violation of the client threshold occurs only if one or more source conditions are violated (zero false negative), and (2) the number of source violations when client threshold is not violated is small (minimize false positives). Using performance evaluation we show that our algorithms result in up to *an order of magnitude* less number of false positives compared to the approaches in the literature.

Probabilistic Top-*k* **Query Processing in Distributed Sensor Networks** *(Short paper)*

Mao Ye¹, Xingjie Liu¹, Wang-Chien Lee¹, Dik Lun Lee²; ¹Pennsylvania State University, USA; ²Hong Kong University of Science & Technology, China

In this paper, we propose the notion of *sufficient set* for distributed processing of probabilistic Top-*k* queries in cluster-based wireless sensor networks. Through the derivation

of sufficient boundary, we show that data items ranked lower than sufficient boundary are not required for answering the probabilistic top-*k* queries, thus are subject to local pruning. Accordingly, we develop the *sufficient set-based (SSB)* algorithm for inter-cluster query processing. Experimental results show that the proposed algorithm reduces data transmissions significantly.

Research Session 18: Query Optimization

Regency C, 16:00 – 17:30, Wednesday Chair: Jingren Zhou

Polynomial Heuristics for Query Optimization

Nicolas Bruno, César Galindo-Legaria, Milind Joshi; Microsoft, USA

Research on query optimization has traditionally focused on exhaustive enumeration of an exponential number of candidate plans. Alternatively, heuristics for query optimization are restricted in several ways, such as by either focusing on join predicates only, ignoring the availability of indexes, or in general having high-degree polynomial complexity. In this paper we propose a heuristic approach to very efficiently obtain execution plans for complex queries, which takes into account the presence of indexes and goes beyond simple join reordering. We also introduce a realistic workload generator and validate our approach using both synthetic and real data.

Optimized Query Evaluation Using Cooperative Sorts

Yu Cao, Ramadhana Bramandia, Chee-Yong Chan, Kian-Lee Tan; National University of Singapore, Singapore

Many applications require sorting a table over multiple sort orders: generation of multiple reports from a table, evaluation of a complex query that involves multiple instances of a relation, and batch processing of a set of queries. In this paper, we study how multiple sortings of a table can be efficiently performed. We introduce a new evaluation technique, called cooperative sort, that exploits the relationships among the input set of sort orders to minimize I/O operations for the collection of sort operations. To demonstrate the efficiency of the proposed scheme, we implemented it in PostgreSQL and evaluated its performance using both TPC-DS benchmark and synthetic data. Our experimental results show significant performance improvement over the traditional non-cooperative sorting scheme.

Generating Code for Holistic Query Evaluation

Konstantinos Krikellas, Stratis D. Viglas, Marcelo Cintra; University of Edinburgh, UK

We present the application of customized code generation to database query evaluation. The idea is to use a collection of highly efficient code templates and dynamically instantiate them to create query- and hardware-specific source code. The source code is compiled and dynamically linked to the database server for processing. Code generation diminishes the bloat of higher-level programming abstractions necessary for implementing generic, interpreted, SQL query engines. At the same time, the generated code is customized for the hardware it will run on. We term this approach *holistic query evaluation*. We present the design and development of a prototype system called HIQUE, the *Holistic Integrated Query Engine*, which incorporates our proposals. We undertake a detailed experimental study of the system's performance. The results show that HIQUE satisfies its design objectives, while its efficiency surpasses that of both well-established and currently-emerging query processing techniques.

Research Session 19: Graph Mining

Regency D, 16:00 – 17:30, Wednesday Chair: Wook-Shin Han

Finding Clusters in Subspaces of Very Large, Multi-Dimensional Datasets

*Robson L.F. Cordeiro*¹, *Agma J.M. Traina*¹, *Christos Faloutsos*², *Caetano Traina Jr.*¹; ¹Universidade de São Paulo, Brazil; ²Carnegie Mellon University, USA

We propose the Multi-resolution Correlation Cluster detection (MrCC), a novel, scalable method to detect correlation clusters able to analyze dimensional data in the range of around 5 to 30 axes. Existing methods typically exhibit super-linear behavior in terms of space or execution time. MrCC employs a novel data structure based on multi-resolution and gains over previous approaches in: (a) it finds clusters that stand out in the data in a statistical sense; (b) it is linear on running time and memory usage regarding number of data points and dimensionality of subspaces where clusters exist; (c) it is linear in memory usage and quasi-linear in running time regarding space dimensionality; and (d) it is accurate, deterministic, robust to noise, does not require stating the number of clusters as input parameter, does not perform distance calculation and is able to detect clusters in subspaces generated by original axes or linear combinations of original axes, including space rotation. We performed experiments on synthetic data ranging from 5 to 30 axes and from 12k to 250k points, and MrCC outperformed in time five of the recent and related work, being in average 10 times faster than the competitors that also presented high accuracy results for every tested dataset. Regarding real data, MrCC found clusters at least 9 times faster than the competitors, increasing their accuracy in up to 34 percent.

Similarity Search on Supergraph Containment

*Haichuan Shang*¹, *Ke Zhu*¹, *Xuemin Lin*¹, *Ying Zhang*¹, *Ryutaro Ichise*²; ¹*University of New South Wales, Australia;* ²*National Institute of Informatics, Japan*

A supergraph containment search is to retrieve the data graphs contained by a query graph. In this paper, we study the problem of efficiently retrieving all data graphs approximately contained by a query graph, namely similarity search on supergraph containment. We propose a novel and efficient index to boost the efficiency of query processing. We have studied the query processing cost and propose two index construction strategies aimed at optimizing the performance of different types of data graphs: top-down strategy and bottom-up strategy. Moreover, a novel indexing technique is proposed by effectively merging the indexes of individual data graphs; this not only reduces the index size but also further reduces the query processing time. We conduct extensive experiments on real data sets to demonstrate the efficiency and the effectiveness of our techniques.

Finding Top-k Maximal Cliques in an Uncertain Graph (Short paper)

Zhaonian Zou, Jianzhong Li, Hong Gao, Shuo Zhang; Harbin Institute of Technology, China

Existing studies on graph mining focus on exact graphs that are precise and complete. However, graph data tends to be uncertain in practice due to noise, incompleteness and inaccuracy. This paper investigates the problem of finding top-*k* maximal cliques in an uncertain graph. A new model of uncertain graphs is presented, and an intuitive measure is introduced to evaluate the significance of vertex sets. An optimized branch-and-bound algorithm is developed to find top-*k* maximal cliques, which adopts efficient pruning rules, a new searching strategy and effective preprocessing methods. The extensive experimental results show that the proposed algorithm is very efficient on real uncertain graphs, and the top-*k* maximal cliques are very useful for real applications, e.g. protein complex prediction.
Progressive Clustering of Networks Using Structure-Connected Order of Traversal (Short paper)

Dustin Bortner, Jiawei Han; University of Illinois at Urbana-Champaign, USA

Network clustering enables us to view a complex network at the macro level, by grouping its nodes into units whose characteristics and interrelationships are easier to analyze and understand. State-of-the-art network partitioning methods are unable to identify hubs and outliers. A recently proposed algorithm, SCAN, overcomes this difficulty. However, it requires a minimum similarity parameter ε but provides no automated way to find it. Thus, it must be rerun for each ε value and does not capture the variety or hierarchy of clusters. We propose a new algorithm, SCOT (or Structure-Connected Order of Traversal), that produces a length *n* sequence containing all possible ε -clusterings. We propose a new algorithm, HintClus (or Hierarchy-Induced Network Clustering), to hierarchically cluster the network by finding only best cluster boundaries (not agglomerative). Results on model-based synthetic network data and real data show that SCOT's execution time is comparable to SCAN, that HintClus runs in negligible time, and that HintClus produces sensible clusters in the presence of noise.

Research Session 20: Parallel Processing

Beacon A, 16:00 – 17:30, Wednesday Chair: Mirek Riedewald

Osprey: Implementing MapReduce-Style Fault Tolerance in a Shared-Nothing Distributed Database

Christopher Yang, Christine Yen, Ceryen Tan, Samuel R. Madden; MIT, USA

In this paper, we describe a scheme for tolerating and recovering from mid-query faults in a distributed shared nothing database. Rather than aborting and restarting queries, our system, *Osprey*, divides running queries into *subqueries*, and replicates data such that each subquery can be rerun on a different node if the node initially responsible fails or returns too slowly. Our approach is inspired by the fault tolerance properties of MapReduce, in which map or reduce jobs are greedily assigned to workers, and failed jobs are rerun on other workers.

Osprey is implemented using a middleware approach, with only a small amount of custom code to handle cluster coordination. Each node in the system is a discrete database system running on a separate machine. Data, in the form of tables, is partitioned amongst database nodes and each partition is replicated on several nodes, using a technique called *chained declustering* [1]. A coordinator machine acts as a standard SQL interface to users; it transforms an input SQL query into a set of subqueries that are then executed on the nodes. Each subquery represents only a small fraction of the total execution of the query; worker nodes are assigned a new subquery as they finish their current one. In this greedy-approach, the amount of work lost due to node failure is small (at most one subquery's work), and the system is automatically load balanced, because slow nodes will be assigned fewer subqueries.

We demonstrate Osprey's viability as a distributed system for a small data warehouse data set and workload. Our experiments show that the overhead introduced by the middleware is small compared to the workload, and that the system shows promising load balancing and fault tolerance properties.

FPGA Acceleration for the Frequent Item Problem

Jens Teubner, Rene Mueller, Gustavo Alonso; ETH Zürich, Switzerland

Field-programmable gate arrays (FPGAs) can provide performance advantages with a lower resource consumption (*e.g.*, energy) than conventional CPUs. In this paper, we show how to employ FPGAs to provide an efficient and high-performance solution for the *frequent item* problem.

We discuss three design alternatives, each one of them exploiting different FPGA features, and we provide an exhaustive evaluation of their performance characteristics. The first design is a one-to-one mapping of the *Space-Saving* algorithm (shown to be the best approach in software [1]), built on special features of FPGAs: *content-addressable memory* and *dualported BRAM*. The two other implementations exploit the flexibility of digital circuits to implement *parallel lookups* and *pipelining strategies*, resulting in significant improvements in performance.

On low-cost FPGA hardware, the fastest of our designs can process 80 million items per second — three times as much as the best known result. Moreover, and unlike in software approaches where performance is directly related to the skew factor of the Zipf distribution, the high throughput is independent of the skew of the distribution of the input. In the paper we discuss as well several design trade-offs that are relevant when implementing database functionality on FPGAs. In particular, we look at *resource consumption* and the levels of *data* and *task parallelism* of three different designs.

Estimating the Progress of MapReduce Pipelines (Short paper)

Kristi Morton, Abram Friesen, Magdalena Balazinska, Dan Grossman; University of Washington, USA

In parallel query-processing environments, *accurate, time-oriented* progress indicators could provide much utility given that inter- and intra-query execution times can have high variance. However, none of the techniques used by existing tools or available in the literature provide non-trivial progress estimation for parallel queries. In this paper, we introduce Parallax, the first such indicator. While several parallel data processing systems exist, the work in this paper targets environments where queries consist of a series of MapReduce jobs. Parallax builds on recently-developed techniques for estimating the progress of single-site SQL queries, but focuses on the challenges related to parallelism and variable execution speeds. We have implemented our estimator in the Pig system and demonstrate its performance through experiments with the PigMix benchmark and other queries running in a real, smallscale cluster.

Scalable Distributed-Memory External Sorting (Short paper)

Mirko Rahn, Peter Sanders, Johannes Singler; Karlsruhe Institute of Technology, Germany

We engineer algorithms for sorting huge data sets on massively parallel machines. The algorithms are based on the multiway merging paradigm. We first outline an algorithm whose I/O requirement is close to a lower bound. Thus, in contrast to naive implementations of multiway merging and all other approaches known to us, the algorithm works with just two passes over the data even for the largest conceivable inputs. A second algorithm reduces communication overhead and uses more conventional specifications of the result at the cost of slightly increased I/O requirements. An implementation wins the well known sorting benchmark in several categories and by a large margin over its competitors.

Research Session 21: Keyword Search

Regency B, 10:30 – 12:00, Thursday Chair: Zhen Liu

Supporting Top-K Keyword Search in XML Databases

Liang Jeff Chen, Yannis Papakonstantinou; University of California at San Diego, USA

Keyword search is considered to be an effective information discovery method for both structured and semi-structured data. In XML keyword search, query semantics is based on the concept of *Lowest Common Ancestor* (LCA). However, naive LCA-based semantics leads to exponential computation and result size. In the literature, LCA-based semantic variants (e.g., ELCA and SLCA) were proposed, which define a subset of all the LCAs as the results. While most existing work focuses on algorithmic efficiency, top-K processing for

XML keyword search is an important issue that has received very little attention. Existing algorithms focusing on efficiency are designed to optimize the semantic pruning and are incapable of supporting top-K processing. On the other hand, straightforward applications of top-K techniques from other areas (e.g., relational databases) generate LCAs that may not be the results and unnecessarily expand efforts in the semantic pruning. In this paper, we propose a series of join-based algorithms that combine the semantic pruning and the top-K processing to support top-K keyword search in XML databases. The algorithms essentially reduce the keyword query evaluation to relational joins, and incorporate the idea of the top-K join from relational databases. Extensive experimental evaluations show the performance advantages of our algorithms.

Personalized Web Search with Location Preferences

*Kenneth Wai-Ting Leung*¹, *Dik Lun Lee*¹, *Wang-Chien Lee*²; ¹*Hong Kong University of Science & Technology, China*; ²*Pennsylvania State University, USA*

As the amount of Web information grows rapidly, search engines must be able to retrieve information according to the user's preference. In this paper, we propose a new web search personalization approach that captures the user's interests and preferences in the form of concepts by mining search results and their clickthroughs. Due to the important role location information plays in mobile search, we separate concepts into content concepts and location concepts, and organize them into ontologies to create an *ontology-based, multi*facet (OMF) profile to precisely capture the user's content and location interests and hence improve the search accuracy. Moreover, recognizing the fact that different users and queries may have different emphases on content and location information, we introduce the notion of content and location entropies to measure the amount of content and location information associated with a query, and *click* content and location entropies to measure how much the user is interested in the content and location information in the results. Accordingly, we propose to define personalization effectiveness based on the entropies and use it to balance the weights between the content and location facets. Finally, based on the derived ontologies and personalization effectiveness, we train an SVM to adapt a personalized ranking function for re-ranking of future search. We conduct extensive experiments to compare the precision produced by our OMF profiles and that of a baseline method. Experimental results show that OMF improves the precision significantly compared to the baseline.

Fuzzy Matching of Web Queries to Structured Data (Short paper)

*Tao Cheng*¹, *Hady W. Lauw*², *Stelios Paparizos*²; ¹*University of Illinois at Urbana-Champaign, USA*; ²*Microsoft, USA*

Recognizing the alternative ways people use to reference an entity, is important for many Web applications that query structured data. In such applications, there is often a mismatch between how content creators describe entities and how different users try to retrieve them. In this paper, we consider the problem of determining whether a candidate query approximately matches with an entity. We propose an off-line, data-driven, bottom-up approach that mines query logs for instances where Web content creators and Web users apply a variety of strings to refer to the same Web pages. This way, given a set of strings that reference entities, we generate an expanded set of equivalent strings for each entity. The proposed method is verified with experiments on real-life data sets showing that we can dramatically increase the queries that can be matched.

Toward Industrial-Strength Keyword Search Systems Over Relational Data (*Short paper*)

Akanksha Baid, Ian Rae, AnHai Doan, Jeffrey F. Naughton; University of Wisconsin-Madison, USA

Keyword search (KWS) over relational data, where the answers are multiple tuples connected via joins, has received significant attention in the past decade. Numerous solutions have been proposed and many prototypes have been developed. Building on this rapid progress and on growing user needs, recently several RDBMS and Web companies as well as academic research groups have started to examine how to build industrial-strength keywords search

systems. This task clearly requires addressing many issues, including robustness, accuracy, reliability, and privacy, among others. A major emerging issue, however, appears to be performance related: current KWS systems have unpredictable run time. In particular, for certain queries it takes too long to produce answers, and for others the system may even fail to return (e.g., after exhausting memory).

In this paper we begin by examining the above problem and arguing that it is a fundamental problem unlikely to be solved in the near future by software and hardware advances. Next, we argue that in an industrial-strength setting, to ensure real-time interaction and facilitate user adoption, KWS systems should produce answers under an absolute time limit and then provide users with a description of what could be done next, should he or she choose to continue. Next, we show how to realize these requirements for DISCOVER, an exemplar of a recent KWS solution approach. Our basic idea is to produce answers as in today's KWS systems up to the time limit, then show users these answers as well as query forms that characterize the unexplored portion of the answer space. Finally, we present some pre-liminary experiments over real-world data to demonstrate the feasibility of the proposed solution approach.

Research Session 22: Query Processing

Regency C, 10:30 – 12:00, Thursday Chair: Florian Waas

Efficient Processing of Substring Match Queries with Inverted q-Gram Indexes

*Younghoon Kim*¹, *Kyoung-Gu Woo*², *Hyoungmin Park*¹, *Kyuseok Shim*¹; ¹*Seoul National University, Korea;* ²*Samsung Electronics, Korea*

With the widespread of the internet, text-based data sources have become ubiquitous and the demand of effective support for string matching queries becomes ever increasing. The relational query language SQL also supports LIKE clause over string data to handle substring matching queries. Due to popularity of such substring matching queries, there have been a lot of study on designing efficient indexes to support the LIKE clause in SQL. Among them, q-gram based indexes have been studied extensively. However, how to process substring matching queries efficiently with such indexes has received very little attention until recently.

In this paper, we show that the optimal execution of intersecting posting lists of q-grams for substring matching queries should be decided judiciously. Then we present the optimal and approximate algorithms based on cost estimation for substring matching queries. Performance study confirms that our techniques improve query execution time with q-gram indexes significantly compared to the traditional algorithms.

Progressive Result Generation for Multi-Criteria Decision Support Queries

Venkatesh Raghavan, Elke A. Rundensteiner; Worcester Polytechnic Institute, USA

Multi-criteria decision support (MCDS) is crucial in many business and web applications such as web searches, B2B portals and on-line commerce. Such MCDS applications need to report results early; as soon as they are being generated so that they can react and formulate competitive decisions in near real-time. The ease in expressing user preferences in web-based applications has made Pareto-optimal (*skyline*) queries a popular class of MCDS queries. However, state-of-the-art techniques either focus on handling skylines on single input sets (i.e., no joins) or do not tackle the challenge of producing progressive early output results. In this work, we propose a progressive query evaluation framework *ProgXe* that transforms the execution of queries involving skyline over joins to be *non-blocking*, i.e., to be progressively generating results early and often. In *ProgXe* the query processing (join, mapping and skyline) is conducted at multiple levels of abstraction, thereby exploiting the knowledge gained from both input as well as mapped output spaces. This knowledge enables us to identify and reason about abstract-level relationships to guarantee correctness of early output. It also provides optimization opportunities previously missed by current techniques. To further optimize *ProgXe*, we incorporate an ordering technique that optimizes the rate at which results are reported by translating the optimization of tuple-level processing into a job-sequencing problem. Our experimental study over a wide variety of data sets demonstrates the superiority of our approach over state-of-the-art techniques.

Nb-GCLOCK: A Non-Blocking Buffer Management Based on the Generalized CLOCK

*Makoto Yui*¹, *Jun Miyazaki*², *Shunsuke Uemura*³, *Hayato Yamana*¹; ¹Waseda University. Japan: ²NAIST. Japan: ³Nara Sanavo University. Iavan

In this paper, we propose a non-blocking buffer management scheme based on a lock-free variant of the GCLOCK page replacement algorithm. Concurrent access to the buffer management module is a major factor that prevents database scalability to processors. Therefore, we propose a non-blocking scheme for bufferfix operations that fix buffer frames for requested pages without locks by combining Nb-GCLOCK and a non-blocking hash table. Our experimental results revealed that our scheme can obtain nearly linear scalability to processors up to 64 processors, although the existing locking-based schemes do not scale beyond 16 processors.

Research Session 23: Web and Collaborative Applications

Regency D, 10:30 - 12:00, Thursday Chair: Zack Ives

Effective Automated Object Matching

Diego Zardetto¹, Monica Scannapieco¹, Tiziana Catarci²; ¹Istituto Nazionale di Statistica, Italy; ²Università di Roma "La Sapienza", Italy

Object Matching (OM) is the problem of identifying pairs of data-objects coming from different sources and representing the same real world object. Several methods have been proposed to solve OM problems, but none of them seems to be at the same time fully automated and very effective. In this paper we present a fundamentally new suite of methods that instead possesses both these abilities.

We adopt a statistical approach based on mixture models, which structures an OM process into two consecutive tasks. First, mixture parameters are estimated by fitting the model to observed distance measures between pairs. Then, a probabilistic clustering of the pairs into Matches and Unmatches is obtained by exploiting the fitted model.

In particular, we use a mixture model with component densities belonging to the Beta parametric family and we fit it by means of an original perturbation-like technique. Moreover, we solve the clustering problem according to both Maximum Likelihood and Minimum Cost objectives. To accomplish this task, optimal decision rules fulfilling one-to-one matching constraints are searched by a purposefully designed evolutionary algorithm.

Notably, our suite of methods is distance-independent in the sense that it does not rely on any restrictive assumption on the function to be used when comparing data-objects. Even more interestingly, our approach is not confined to record linkage applications but can be applied to match also other kinds of data-objects. We present several experiments on real data that validate the proposed methods and show their excellent effectiveness.

Efficient Identification of Coupled Entities in Document Collections

(Short paper)

Nikos Sarkas¹, Albert Angel¹, Nick Koudas¹, Divesh Srivastava²; ¹University of Toronto, Canada; ²AT&T Labs Research, USA

The relentless pace at which textual data are generated on-line necessitates novel paradigms for their understanding and exploration. To this end, we introduce a methodology for discovering *strong entity associations* in all the *slices* (meta-data value restrictions) of a document collection. Since related documents mention approximately the same group of core entities (people, locations, etc.), the groups of coupled entities discovered can be used to expose themes in the document collection.

We devise and evaluate algorithms capable of addressing two flavors of our core problem: algorithm THR-ENT for computing all sufficiently strong entity associations and algorithm TOP-ENT for computing the top-k strongest entity associations, for each slice of the document collection.

On Supporting Effective Web Extraction (Short paper)

*Wook-Shin Han*¹, *Wooseong Kwak*¹, *Hwanjo Yu*²; ¹*Kyungpook National University, Korea*; ²*POSTECH, Korea*

Commercial tuple extraction systems have enjoyed some success to extract tuples by regarding HTML pages as tree structures and exploiting XPath queries to find attributes of tuples in the HTML pages. However, such systems would be vulnerable to small changes on the web pages. In this paper, we propose a robust tuple extraction system which utilizes spatial relationships among elements rather than the XPath queries of the elements. Our system regards elements in the rendered page as spatial objects in the 2-D space and executes spatial joins to extract target elements. Since humans also identify an element in a web page by its relative spatial location, our system extracting elements by their spatial relationships could possibly be as robust as manual extraction and is far more robust than existing tuple extraction systems.

A Partial Persistent Data Structure to Support Consistency in Real-Time Collaborative Editing (Short paper)

Qinyi Wu¹, Calton Pu¹, João Eduardo Ferreiar²; ¹Georgia Institute of Technology, USA; ²Universidade de São Paulo, Brazil

Co-authored documents are becoming increasingly important for knowledge representation and sharing. Tools for supporting document co-authoring are expected to satisfy two requirements: 1) querying changes over editing histories; 2) maintaining data consistency among users. Current tools support either limited queries or are not suitable for loosely controlled collaborative editing scenarios. We address both problems by proposing a new persistent data structure — partial persistent sequence. The new data structure enables us to create unique character identifiers that can be used for associating meta-information and tracking their changes, and also design simple view synchronization algorithms to guarantee data consistency under the presence of concurrent updates. Experiments based on real-world collaborative editing traces show that our data structure uses disk space economically and provides efficient performance for document update and retrieval.

Detecting Bursty Events in Collaborative Tagging Systems (Short paper)

Junjie Yao, Bin Cui, Yuxin Huang, Yanhong Zhou; Peking University, China

Collaborative tagging systems have emerged as an ubiquitous way to annotate and organize online resources. The users' tagging actions over time reflect the changing of their interests. In this paper, we propose to detect *bursty tagging event*, which captures the relations among a group of correlated tags where the tags are either bursty or associated with bursty tag co-occurrence. We exploit the sliding time intervals to extract bursty features from large tag corpora as the first step, and then adopt graph clustering techniques to group bursty features into meaningful bursty events. An experimental study demonstrates the superiority of our approach.

Research Session 24: Scientific Databases

Regency B, 16:00 – 17:30, Thursday Chair: Feifei Li

Credibility-Enhanced Curated Database: Improving the Value of Curated Databases

Qun Ni, Elisa Bertino; Purdue University, USA

In curated databases, annotations may contain opinions different from those in sources. Moreover, annotations may contradict each other and have uncertainty. Such situations result in a natural question: "Which opinion is most likely to be correct?" In this paper, we define a credibility-enhanced curated database and propose an efficient method to accurately evaluate the correctness of sources and annotations in curated databases.

UV-Diagram: A Voronoi Diagram for Uncertain Data

*Reynold Cheng*¹, *Xike Xie*¹, *Man Lung Yiu*², *Jinchuan Chen*³, *Liwen Sun*¹; ¹*University of Hong Kong, China;* ²*Hong Kong Polytechnic University, China;* ³*Renmin University of China, China*

The Voronoi diagram is an important technique for answering nearest-neighbor queries for spatial databases. In this paper, we study how the Voronoi diagram can be used on uncertain data, which are inherent in scientific and business applications. In particular, we propose the *Uncertain-Voronoi Diagram* (or *UV-diagram* in short). Conceptually, the data space is divided into distinct "UV-partitions", where each UV-partition P is associated with a set S of objects; any point q located in P has the set S as its nearest neighbor with non-zero probabilities. The UV-diagram facilitates queries that inquire objects for having non-zero chances of being the nearest neighbor of a given query point. It also allows analysis of nearest neighbor information, e.g., finding out how many objects are the nearest neighbors in a given area.

However, a UV-diagram requires exponential construction and storage costs. To tackle these problems, we devise an alternative representation for UV-partitions, and develop an adaptive index for the UV-diagram. This index can be constructed in polynomial time. We examine how it can be extended to support other related queries. We also perform extensive experiments to validate the effectiveness of our approach.

Supporting Real-World Activities in Database Management Systems *(Short paper)*

Mohamed Y. Eltabakh, Walid G. Aref, Ahmed K. Elmagarmid, Yasin N. Silva, Mourad Ouzzani; Purdue University, USA

The cycle of processing the data in many application domains is complex and may involve real-world activities that are external to the database, e.g., wet-lab experiments, instrument readings, and manual measurements. These real-world activities may take long time to prepare for and to perform, and hence introduce inherently long time delays between the updates in the database. The presence of these long delays between the updates, along with the need for the intermediate results to be instantly available, makes supporting real-world activities in the database engine a challenging task. In this paper, we address these challenges through a system that enables users to reflect their updates immediately into the database while keeping track of the dependent and *potentially invalid* data items until they are re-validated. The proposed system includes: (1) semantics and syntax for interfaces through which users can express the dependencies among data items, (2) new operators to alert users when the returned query results contain potentially invalid or out-of-date data, and to enable evaluating queries on either valid data only, or both valid and potentially invalid data, and (3) mechanisms for data invalidation and revalidation. The proposed system is being realized via extensions to PostgreSQL.

XML-Based Computation for Scientific Workflows (Short paper)

Daniel Zinn¹, Shawn Bowers², Bertram Ludäscher¹; ¹University of California at Davis, USA; ²Gonzaga University, USA

Scientific workflows are increasingly used to rapidly integrate existing algorithms to create larger and more complex programs. However, designing workflows using purely datafloworiented computation models introduces a number of challenges, including the need to use low-level components to mediate and transform data (so-called *shims*) and large numbers of additional "wires" for routing data to components within a workflow. To address these problems, we employ *Virtual Data Assembly Lines (VDAL)*, a modeling paradigm that can eliminate most shims and reduce wiring complexity. We show how a VDAL design can be implemented using existing XML technologies and how static analysis can provide significant help to scientists during workflow design and evolution, e.g., by displaying actor dependencies or by detecting so-called unproductive actors.

Research Session 25: Tree Queries and Semi-Structured Databases

Regency C, 16:00 – 17:30, Thursday Chair: Maurice van Keulen

ViewJoin: Efficient View-Based Evaluation of Tree Pattern Queries

Ding Chen, Chee-Yong Chan; National University of Singapore, Singapore

There is a lot of recent interest in applying views to optimize the processing of tree pattern queries (TPQs). However, existing work in this area has focused predominantly on logical optimization issues, namely, view selection and query rewriting. With the exception of the recent work on InterJoin (which is primarily focused on path queries and views), there is very little work that has examined the important physical optimization issue of how to efficiently evaluate TPQs using materialized views. In this paper, we present a new storage scheme for materialized TPQ views and a novel evaluation algorithm for processing general TPQ queries using materialized TPQ views. Our experimental results demonstrate that our proposed method outperforms the state-of-the-art approaches.

FlexPref: A Framework for Extensible Preference Evaluation in Database Systems

Justin J. Levandoski, Mohamed F. Mokbel, Mohamed E. Khalefa; University of Minnesota, USA

Personalized database systems give users answers tailored to their personal preferences. While numerous preference evaluation methods for databases have been proposed (e.g., skyline, top-k, k-dominance, k-frequency), the implementation of these methods at the core of a database system is a double-edged sword. Core implementation provides efficient query processing for arbitrary database queries, however this approach is not practical as *each* existing (and future) preference method requires a custom query processor implementation. To solve this problem, this paper introduces FlexPref, a framework for extensible preference evaluation in database systems. FlexPref, implemented in the query processor, aims to support a wide-array of preference evaluation methods in a single extensible code base. Integration with FlexPref is simple, involving the registration of only *three* functions that capture the essence of the preference method. Once integrated, the preference method "lives" at the core of the database, enabling the efficient execution of preference queries involving common database operations. To demonstrate the extensibility of FlexPref, we provide case studies showing the implementation of three database operations (single table access, join, and sorted list access) and five state-of-the-art preference evaluation methods (top-k, skyline, k-dominance, top-k dominance, and k-frequency). We also experimentally study the strengths and weaknesses of an implementation of FlexPef in PostgreSOL over a range of single-table and multi-table preference queries.

Optimal Tree Node Ordering for Child/Descendant Navigations (Short paper)

Atsuyuki Morishima¹, Keishi Tajima², Masateru Tadaishi¹; ¹University of Tsukuba, Japan; ²Kyoto University, Japan

There are many applications in which users interactively access huge tree data by repeating set-based navigations. In this paper, we focus on label-specific/wildcard children/ descendant navigations. For efficient processing of these operations in huge data stored on a disk, we need a node ordering scheme that clusters nodes that are accessed together by these operations. In this paper, (1) we show there is no node order that is optimal for all these operations, (2) we propose two schemes, each of which is optimal only for some subset of them, and (3) we show that one of the proposed schemes can process all these operations with access to a constant-bounded number of regions on the disk without accessing irrelevant nodes.

XMorph: A Shape-Polymorphic, Domain-Specific XML Data Transformation Language (Short paper)

*Curtis Dyreson*¹, *Sourav Bhowmick*², *Aswani Rao Jannu*¹, *Kirankanth Mallampalli*¹, *Shuohao Zhang*³; ¹Utah State University, USA; ²Nanyang Technological University, Singapore; ³Marvel, USA

By imposing a single hierarchy on data, XML makes queries *brittle* in the sense that a query might fail to produce the desired result if it is executed on the same data organized in a different hierarchy, or if the hierarchy evolves during the lifetime of an application. This paper presents a new transformation language, called XMorph, which supports more flexible querying. XMorph is a shape polymorphic language, that is, a single XMorph query can extract and transform data from differently-shaped hierarchies. The XMorph data shredder distills XML data into a graph of *closest relationships*, which are exploited by the query evaluation engine to produce a result in the shape specified by an XMorph query.

Research Session 26: Query Ranking and Database Testing

Regency D, 16:00 – 17:30, Thursday Chair: Cesar Galindo-Legaria

Surrogate Ranking for Very Expensive Similarity Queries

*Fei Xu*¹, *Ravi Jampani*¹, *Mingxi Wu*², *Chris Jermaine*¹, *Tamer Kahveci*¹; ¹*University of Florida, USA*; ²*Oracle, USA*

We consider the problem of similarity search in applications where the cost of computing the similarity between two records is very expensive, and the similarity measure is not a metric. In such applications, comparing even a tiny fraction of the database records to a single query record can be orders of magnitude slower than reading the entire database from disk, and indexing is often not possible. We develop a general-purpose, statistical framework for answering top-k queries in such databases, when the database administrator is able to supply an inexpensive surrogate ranking function that substitutes for the actual similarity measure. We develop a robust method that learns the relationship between the surrogate function and the similarity measure. Given a query, we use Bayesian statistics to update the model by taking into account the observed partial results. Using the updated model, we construct bounds on the accuracy of the result set obtained via the surrogate ranking. Our experiments show that our models can produce useful bounds for several real-life applications.

Semantic Ranking and Result Visualization for Life Sciences Publications

Julia Stoyanovich, William Mee, Kenneth A. Ross; Columbia University, USA An ever-increasing amount of data and semantic knowledge in the domain of life sciences is bringing about new data management challenges. In this paper we focus on adding the semantic dimension to literature search, a central task in scientific research. We focus our attention on PubMed, the most significant bibliographic source in life sciences, and explore ways to use high-quality semantic annotations from the MeSH vocabulary to rank search results. We start by developing several families of ranking functions that relate a search query to a document's annotations. We then propose an efficient adaptive ranking mechanism for each of the families. We also describe a two-dimensional Skyline-based visualization that can be used in conjunction with the ranking to further improve the user's interaction with the system, and demonstrate how such Skylines can be computed adaptively and efficiently. Finally, we evaluate the effectiveness of our ranking with a user study.

Ranked Queries Over Sources with Boolean Query Interfaces without Ranking Support (Short paper)

*Vagelis Hristidis*¹, *Yuheng Hu*¹, *Panagiotis G. Ipeirotis*²; ¹*Florida International University, USA*; ²*New York University, USA*

Many online or local data sources provide powerful querving mechanisms but limited ranking capabilities. For instance, PubMed allows users to submit highly expressive Boolean keyword queries, but ranks the query results by date only. However, a user would typically prefer a ranking by relevance, measured by an Information Retrieval (IR) ranking function. The naive approach would be to submit a disjunctive query with all query keywords, retrieve the returned documents, and then re-rank them. Unfortunately, such an operation would be very expensive due to the large number of results returned by disjunctive queries. In this paper we present algorithms that return the top results for a query, ranked according to an IR-style ranking function, while operating on top of a source with a Boolean query interface with no ranking capabilities (or a ranking capability of no interest to the end user). The algorithms generate a series of conjunctive queries that return only documents that are candidates for being highly ranked according to a relevance metric. Our approach can also be applied to other settings where the ranking is monotonic on a set of factors (query keywords in IR) and the source query interface is a Boolean expression of these factors. Our comprehensive experimental evaluation on the PubMed database and TREC dataset show that we achieve order of magnitude improvement compared to the current baseline approaches.

X-Data: Generating Test Data for Killing SQL Mutants (Short paper)

Bhanu Pratap Gupta, Devang Vira, S. Sudarshan; IIT Bombay, India

Checking if an SQL query has been written correctly is not an easy task. Formal verification is not applicable, since it is based on comparing a specification with an implementation, whereas SQL queries are essentially a specification without any implementation. Thus, the standard approach for testing queries is to manually check query results on test datasets. Intuitively, a mutant is a query variant that could have been the correct query if the query was in error; a mutant is killed by a dataset if the original query and the mutant return different results on the dataset.

In this paper, we address the problem of generation of test data for an SQL query, to kill mutants. Our work focuses in particular on a class of join/outer-join mutants, which are a common cause of error. To minimize human effort in testing, our techniques generate a test suite containing small and intuitive test datasets, combining them into a single dataset where possible. In the absence of foreign-key constraints, and under certain assumptions, the test suite is complete, i.e. it kills all non-equivalent mutations, in the class of join-type mutations that we consider. We also consider some common types of where-clause predicate mutants. Our techniques have been implemented in a prototype data generation tool.

Research Session 27: Social Networks and Similarity Queries

Regency B, 08:30 – 10:00, Friday Chair: K. Selcuk Candan

Discovery-Driven Graph Summarization

*Ning Zhang*¹, *Yuanyuan Tian*², *Jignesh M. Patel*¹; ¹*University of Wisconsin-Madison, USA*; ²*IBM, USA*

Large graph datasets are ubiquitous in many domains, including social networking and biology. Graph summarization techniques are crucial in such domains as they can assist in uncovering useful insights about the patterns hidden in the underlying data. One important type of graph summarization is to produce small and informative summaries based on user-selected node attributes and relationships, and allowing users to interactively drilldown or roll-up to navigate through summaries with different resolutions. However, two key components are missing from the previous work in this area that limit the use of this method in practice. First, the previous work only deals with categorical node attributes. Consequently, users have to manually bucketize numerical attributes based on domain knowledge, which is not always possible. Moreover, users often have to manually iterate through many resolutions of summaries to identify the most interesting ones. This paper addresses both these key issues to make the interactive graph summarization approach more useful in practice. We first present a method to automatically categorize numerical attributes values by exploiting the domain knowledge hidden inside the node attributes values and graph link structures. Furthermore, we propose an interestingness measure for graph summaries to point users to the potentially most insightful summaries. Using two real datasets, we demonstrate the effectiveness and efficiency of our techniques.

The Similarity Join Database Operator

*Yasin N. Silva*¹, *Walid G. Aref*¹, *Mohamed H. Ali*²; ¹*Purdue University, USA*; ²*Microsoft, USA*

Similarity joins have been studied as key operations in multiple application domains, e.g., record linkage, data cleaning, multimedia and video applications, and phenomena detection on sensor networks. Multiple similarity join algorithms and implementation techniques have been proposed. They range from out-of-database approaches for only in-memory and external memory data to techniques that make use of standard database operators to answer similarity joins. Unfortunately, there has not been much study on the role and implementation of similarity joins as database physical operators. In this paper, we focus on the study of similarity joins as first-class database operators. We present the definition of several similarity join operators and study the way they interact among themselves, with other standard database operators, and with other previously proposed similarity-aware operators. In particular, we present multiple transformation rules that enable similarity query optimization through the generation of equivalent similarity query execution plans. We then describe an efficient implementation of two similarity join operators, ϵ -Join and Join-Around, as core DBMS operators. The performance evaluation of the implemented operators in PostgreSQL shows that they have good execution time and scalability properties. The execution time of Join-Around is less than 5% of the one of the equivalent query that uses only regular operators while ε -Join's execution time is 20 to 90% of the one of its equivalent regular operators based query for the useful case of small ε (0.01% to 10% of the domain range). We also show experimentally that the proposed transformation rules can generate plans with execution times that are only 10% to 70% of the ones of the initial query plans.

Anonymizing Weighted Social Network Graphs (Short paper)

Sudipto Das, Ömer Eğecioğlu, Amr El Abbadi; University of California at Santa Barbara, USA

The increasing popularity of social networks has initiated a fertile research area in infor-

mation extraction and data mining. Although such analysis can facilitate better understanding of sociological, behavioral, and other interesting phenomena, there is a growing concern about personal privacy being breached, thereby requiring effective anonymization techniques. In this paper, we consider edge weight anonymization in social graphs. Our approach builds a linear programming (LP) model which preserves properties of the graph that are expressible as linear functions of the edge weights. Such properties form the foundations of many important graph-theoretic algorithms such as *shortest paths, k-nearest neighbors, minimum spanning tree*, etc. Off-the-shelf LP solvers can then be used to find solutions to the resulting model where the computed solution constitutes the weights in the anonymized graph. As a proof of concept, we choose the *shortest paths problem*, and experimentally evaluate the proposed techniques using real social network data sets.

Efficient Similarity Matching of Time Series Cliques with Natural Relations (Short paper)

*Zhe Zhao*¹, *Bin Cui*¹, *Wee Hyong Tok*², *Jiakui Zhao*³; ¹*Peking University, China*; ²*Microsoft, China*; ³*China Electric Power Research Institute, China*

A Time Series Clique (*TSC*) consists of multiple time series. In each TSC, the time series hold some natural relations with each other. In conventional time series retrieval methods, such natural relations are often ignored. In this paper, we formalize the problem of similarity search over TSC databases and develop a novel framework for similarity search on TSC data, which considers both time series patterns and relations. We conduct an extensive performance study, and the results show the effectiveness and efficiency of the proposed method.

Research Session 28: Stream Processing

Regency C, 08:30 – 10:00, Friday Chair: Alexandros Labrinidis

Continuous Query Evaluation Over Distributed Sensor Networks

Oana Jurca, Sebastian Michel, Alexandre Herrmann, Karl Aberer; EPFL, Switzerland

In this paper we address the problem of processing continuous multi-join queries, over distributed data streams. Our approach makes use of existing work in the field of publish/subscribe systems. We show how these principles can be ported to our envisioned architectural model by enriching the common query model with location dependent attributes. We allow users to subscribe to a set of sensor attributes, a service that requires processing multi-join correlation queries. The goal is to decrease the overall network traffic consumption by removing redundant subscriptions and eliminating unrequested events close to the publishing sensors. This is non-trivial, especially in the presence of multi-join queries without any central control mechanism. Our approach is based on the concept of filter-split-forward phases for efficient subscription filtering and placement inside the network. We report on a performance evaluation using a real-world dataset, showing the improvements over the state-of-the-art, as we reduce the overall data traffic by half.

Space-Efficient Online Approximation of Time Series Data: Streams, Amnesia, and Out-of-Order

Sorabh Gandhi, Luca Foschini, Subhash Suri; University of California at Santa Barbara, USA

In this paper, we present an abstract framework for online approximation of time-series data that yields a unified set of algorithms for several popular models: data streams, amnesic approximation, and out-of-order stream approximation. Our framework essentially develops a popular greedy method of bucket-merging into a more generic form, for which we can prove space-quality approximation bounds. When specialized to piecewise linear bucket approximations and commonly used error metrics, such as L_2 or L_{∞} , our framework leads to provable error bounds where none were known before, offers new results, or yields

simpler and unified algorithms. The conceptual simplicity of our scheme translates into highly practical implementations, as borne out in our simulation studies: the algorithms produce near-optimal approximations, require very small memory footprints, and run extremely fast.

Approximation Trade-Offs in Markovian Stream Processing: An Empirical Study (Short paper)

*Julie Letchner*¹, *Christopher Ré*², *Magdalena Balazinska*¹, *Matthai Philipose*³; ¹University of Washington, USA; ²University of Wisconsin-Madison, USA; ³Intel, USA

A large amount of the world's data is both *sequential* and *imprecise*. Such data is commonly modeled as *Markovian streams*; examples include words/sentences inferred from raw audio signals, or discrete location sequences inferred from RFID or GPS data. The rich semantics and large volumes of these streams make them difficult to query efficiently. In this paper, we study the effects — on both efficiency and accuracy — of two common stream approximations. Through experiments on a real-world RFID data set, we identify conditions under which these approximations can improve performance by several orders of magnitude, with only minimal effects on query results. We also identify cases when the full rich semantics are necessary.

FENCE: Continuous Access Control Enforcement in Dynamic Data Stream Environments (Short paper)

*Rimma V. Nehme*¹, *Hyo-Sang Lim*², *Elisa Bertino*²; ¹*Microsoft, USA*; ²*Purdue University, USA*

In this paper, we present *FENCE* framework that addresses the problem of *continuous access control enforcement* in dynamic data stream environments. The distinguishing characteristics of *FENCE* include: (1) the *stream-centric* approach to security, (2) the *symmetric* modeling of security for both continuous queries and streaming data, and (3) *security-aware query processing* that considers both regular and security-related selectivities. In *FENCE*, both data and query security restrictions are modeled in the form of streaming security metadata, called "*security punctuations*", embedded inside data streams. We have implemented *FENCE* in a prototype DSMS and briefly summarize our performance observations.

Research Session 29: Publishing Privacy

Regency B, 10:30 – 12:00, Friday Chair: Brian Cooper

A Privacy-Preserving Approach to Policy-Based Content Dissemination

Ning Shang, Mohamed Nabeel, Federica Paci, Elisa Bertino; Purdue University, USA

We propose a novel scheme for selective distribution of content, encoded as documents, that preserves the privacy of the users to whom the documents are delivered and is based on an efficient and novel group key management scheme.

Our document broadcasting approach is based on access control policies specifying which users can access which documents, or subdocuments. Based on such policies, a broadcast document is segmented into multiple subdocuments, each encrypted with a different key. In line with modern attribute-based access control, policies are specified against identity attributes of users. However our broadcasting approach is privacy-preserving in that users are granted access to a specific document, or subdocument, according to the policies without the need of providing in clear information about their identity attributes to the document publisher. Under our approach, not only does the document publisher not learn the values of the identity attributes of users, but it also does not learn which policy conditions are verified by which users, thus inferences about the values of identity attributes are prevented. Moreover, our key management scheme on which the proposed broadcasting approach is based is efficient in that it does not require to send the decryption keys to the users along with the encrypted document. Users are able to reconstruct the keys to decrypt the authorized portions of a document based on subscription information they have received from the document publisher. The scheme also efficiently handles new subscription of users and revocation of subscriptions.

Global Privacy Guarantee in Serial Data Publishing (Short paper)

*Raymond Chi-Wing Wong*¹, *Ada Wai-Chee Fu*², *Jia Liu*², *Ke Wang*³, *Yabo Xu*⁴; ¹*Hong Kong University of Science & Technology, China*; ²*Chinese University of Hong Kong, China*; ³*Simon Fraser University, Canada*; ⁴*Sun Yat-sen University, China*

While previous works on privacy-preserving serial data publishing consider the scenario where sensitive values may persist over multiple data releases, we find that no previous work has sufficient protection provided for sensitive values that can change over time, which should be the more common case. In this work, we propose to study the privacy guarantee for such transient sensitive values, which we call the *global guarantee*. We for mally define the problem for achieving this guarantee. We show that the data satisfying the global guarantee also satisfies a privacy guarantee commonly adopted in the privacy literature called the *local guarantee*.

XColor: Protecting General Proximity Privacy (Short paper)

Ting Wang, Ling Liu; Georgia Institute of Technology, USA

As a severe threat in anonymized data publication, proximity breach is gaining increasing attention. Such breach occurs when an attacker learns with high confidence that the sensitive information of a victim associates with a set of semantically proximate values, even though not sure about the exact one. Recently (ε , δ)-dissimilarity [14] has been proposed as an effective countermeasure against general proximity attack. In this paper, we present a detailed analytical study on the fulfillment of this principle, derive criteria to efficiently test its satisfiability for given microdata, and point to a novel anonymization model, XCOLOR, with theoretical guarantees on both operation efficiency and utility preservation.

Correlation Hiding by Independence Masking (Short paper)

Yufei Tao¹, Jian Pei², Jiexing Li¹, Xiaokui Xiao³, Ke Yi⁴, Zhengzheng Xing²; ¹Chinese University of Hong Kong, China; ²Simon Fraser University, Canada; ³Nanyang Technological University, Singapore; ⁴Hong Kong University of Science & Technology, China

Extracting useful correlation from a dataset has been extensively studied. In this paper, we deal with the opposite, namely, a problem we call *correlation hiding* (CH), which is fundamental in numerous applications that need to disseminate data containing sensitive information. In this problem, we are given a relational table T whose attributes can be classified into three disjoint sets A, B, and C. The objective is to distort some values in T so that A becomes independent from B, and yet, their correlation with C is preserved as much as possible. CH is different from all the problems studied previously in the area of data privacy, in that CH demands *complete* elimination of the correlation between two sets of attributes, whereas the previous research focuses on *partial* elimination up to a certain level. A new operator called *independence masking* is proposed to solve the CH problem. Implementations of the operator with good worst case guarantees are described in the full version of this short note.

Regency C, 10:30 – 12:00, Friday Chair: Walid Aref

Monitoring Continuous State Violation in Datacenters: Exploring the Time Dimension

Shicong Meng, Ting Wang, Ling Liu; Georgia Institute of Technology, USA

Monitoring global states of an application deployed over distributed nodes becomes prevalent in today's datacenters. State monitoring requires not only correct monitoring results but also minimum communication cost for efficiency and scalability. Most existing work adopts an instantaneous state monitoring approach, which triggers state alerts whenever a constraint is violated. Such an approach, however, may cause frequent and unnecessary state alerts due to unpredictable monitored value bursts and momentary outliers that are common in large-scale Internet applications. These false alerts may further lead to expensive and problematic counter-measures.

To address this issue, we introduce window-based state monitoring in this paper. Windowbased state monitoring evaluates whether state violation is continuous within a time window, and thus, gains immunity to short-term value bursts and outliers. Furthermore, we find that exploring the monitoring time window at distributed nodes achieves significant communication savings over instantaneous monitoring. Based on this finding, we develop WISE, a system that efficiently performs WIndow-based StatE monitoring at datacenter-scale. WISE is highlighted with three sets of techniques. First, WISE uses distributed filtering time windows and intelligently avoids global information collecting to achieve communication efficiency, while guaranteeing monitoring correctness at the same time. Second, WISE provides a suite of performance tuning techniques to minimize communication cost based on a sophisticated cost model. Third, WISE also employs a set of novel performance optimization techniques. Extensive experiments over both real world and synthetic traces show that WISE achieves a 50%–90% reduction in communication cost compared with existing instantaneous monitoring approaches and simple alternative schemes.

Cost-Efficient and Differentiated Data Availability Guarantees in Data Clouds (Short paper)

Nicolas Bonvin, Thanasis G. Papaioannou, Karl Aberer; EPFL, Switzerland

Failures of any type are common in current data-centers. As data scales up, its availability becomes more complex, while different availability levels per application or per data item may be required. In this paper, we propose a self-managed key-value store that dynamically allocates the resources of a data cloud to several applications in a cost-efficient and fair way. Our approach offers and dynamically maintains multiple differentiated availability guarantees to each different application despite failures. We employ a virtual economy, where each data partition acts as an individual optimizer and chooses whether to migrate, replicate or remove itself based on net benefit maximization regarding the utility offered by the partition and its storage and maintenance cost. Comprehensive experimental evaluation suggest that our solution is highly scalable and adaptive to query rate variations and to resource upgrades/failures.

Intensional Associations in Dataspaces (Short paper)

*Marcos Antonio Vaz Salles*¹, *Jens Dittrich*², *Lukas Blunschi*³; ¹*Cornell University, USA*; ²*Saarland University, Germany*; ³*ETH Zürich, Switzerland*

Dataspace applications necessitate the creation of associations among data items over time. For example, once information about people is extracted from sources on the Web, associations among them may emerge as a consequence of different criteria, such as their city of origin or their elected hobbies. In this paper, we advocate a declarative approach to specifying these associations. We propose that each set of associations be defined by an *association trail*. An association trail is a query-based definition of how items are connected by intensional (i.e., virtual) association edges to other items in the dataspace. We study the problem

of processing neighborhood queries over such intensional association graphs. The naive approach to neighborhood query processing over intensional graphs is to materialize the whole graph and then apply previous work on dataspace graph indexing to answer queries. We present in this paper a novel indexing technique, the grouping-compressed index (GCI), that has better worst-case indexing cost than the naive approach. In our experiments, GCI is shown to provide an order of magnitude gain in indexing cost over the naive approach, while remaining competitive in query processing time.

A Tuple Space for Social Networking on Mobile Phones (Short paper)

Emre Sarigöl, Oriana Riva, Gustavo Alonso; ETH Zürich, Switzerland

Social networking is increasingly becoming a popular means of communication for online users. The trend is also true for offline scenarios where people use their mobile phones to network with nearby buddies. In this paper, we propose a distributed tuple space for social networking on ad hoc networks. We describe the tuple space model and its operations, and give evidence of its advantages for ad hoc social networking through several applications.

Overlapping Community Search for Social Networks (Short paper)

Arnau Padrol-Sureda, Guillem Perarnau-Llobet, Julian Pfeifle, Victor Muntés-Muleor; Universitat Politècnica de Catalunya, Spain

Finding decompositions of a graph into a family of clusters is crucial to understanding its underlying structure. While most existing approaches focus on partitioning the nodes, real-world datasets suggest the presence of overlapping communities. We present OCA, a novel algorithm to detect overlapped communities in large data graphs. It outperforms previous proposals in terms of execution time, and efficiently handles large graphs containing more than 10^8 nodes and edges.

Industry Session 1: Data Warehousing

Regency EF, 10:30 – 12:00, Tuesday Chair: Paul Larson

Hive — A Petabyte Scale Data Warehouse Using Hadoop

Ashish Thusoo, Joydeep Sen Sarma, Namit Jain, Zheng Shao, Prasad Chakka, Ning Zhang, Suresh Antony, Hao Liu, Raghotham Murthy; Facebook, USA

The size of data sets being collected and analyzed in the industry for business intelligence is growing rapidly, making traditional warehousing solutions prohibitively expensive. Hadoop [1] is a popular open-source map-reduce implementation which is being used in companies like Yahoo, Facebook etc. to store and process extremely large data sets on commodity hardware. However, the map-reduce programming model is very low level and requires developers to write custom programs which are hard to maintain and reuse. In this paper, we present Hive, an open-source data warehousing solution built on top of Hadoop. Hive supports queries expressed in a SQL-like declarative language — *HiveQL*, which are compiled into map-reduce jobs that are executed using Hadoop. In addition, HiveQL enables users to plug in custom map-reduce scripts into queries. The language includes a type system with support for tables containing primitive types, collections like arrays and maps, and nested compositions of the same. The underlying IO libraries can be extended to query data in custom formats. Hive also includes a system catalog — *Metastore* — that contains schemas and statistics, which are useful in data exploration, query optimization and query compilation. In Facebook, the Hive warehouse contains tens of thousands of tables and stores over 700TB of data and is being used extensively for both reporting and ad-hoc analyses by more than 200 users per month.

Tuning Servers, Storage and Database for Energy Efficient Data Warehouses

Meikel Poess¹, Raghunath Othayoth Nambiar²; ¹Oracle, USA; ²HP, USA

Undoubtedly, reducing power consumption is at the top of the priority list for system vendors, data center managers who are challenged by customers, analysts, and government agencies to implement green initiatives. Hardware and software vendors have developed an array of power preserving techniques. On-demand-driven clock speeds for processors, energy efficient power supplies, and operating-system-controlled dynamic power modes are just a few hardware examples. Software vendors have contributed to energy efficiency by implementing power efficient coding methods, such as advanced compression and enabling applications to take advantage of large memory caches. However, adoption of these power-preserving technologies in data centers is not straightforward, especially, for large, complex applications such as data warehouses. Data warehouse workloads typically have oscillating resource utilizations, which makes identifying the largest power consumers difficult. Most importantly, while preserving power remains a critical consideration, performance and availability goals must still be met with systems using power-preserving technologies. This paper evaluates the tradeoffs between existing power-saving techniques and their performance impact on data warehouse applications. Our analysis will guide system developers and data center managers in making informed decisions regarding adopting power-preserving techniques.

A New Algorithm for Small-Large Table Outer Joins in Parallel DBMS

Yu Xu, Pekka Kostamaa; Teradata, USA

Large enterprises have been relying on parallel database management systems (*PDBMS*) to process their ever-increasing data volume and complex queries. Business intelligence tools used by enterprises frequently generate a large number of outer joins and require high performance from the underlying database systems. A common type of outer joins in business applications is the small-large table outer join studied in this paper where one table is relatively small and the other is large. We present an efficient and easy to implement algorithm called DER (Duplication and Efficient Redistribution) for small and large table outer joins. Our experimental results show that the DER algorithm significantly speeds up query elapsed time and scales linearly.

Industry Session 2: Data, Data, and More Data

Regency EF, 10:30 – 12:00, Wednesday Chair: Surajit Chaudhuri

Data Cleansing as a Transient Service

Tanveer A. Faruquie, Hima Prasad K., L. Venkata Subramaniam, Mukesh Mohania, Girish Venkatachaliah, Shrinivas Kulkarni, Pramit Basu; IBM, India

There is often a transient need within enterprises for data cleansing which can be satisfied by offering data cleansing as a transient service. Every time a data cleansing need arises it should be possible to provision hardware, software and staff for accomplishing the task and then dismantling the set up. In this paper we present such a system that uses virtualized hardware and software for data cleansing. We share actual experiences gained from building such a system. We use a cloud infrastructure to offer virtualized data cleansing instances that can be accessed as a service. We build a system that is scalable, elastic and configurable. Each enterprise has unique needs which makes it necessary to customize both the infrastructure and the cleansing algorithms to address these needs. In this paper we will present a system that is easily configurable to suit the data cleansing needs of an enterprise.

XBRL Repository — An Industrial approach of Management of XBRL Documents

Zhen Hua Liu, Thomas Baby, Sriram Krishnamurthy, Ying Lu, Qin Yu, Anguel Novoselsky, Vikas Arora; Oracle, USA

XBRL (extensible Business Reporting Language) is an XML based standard for electronic data exchange and communication of business financial documents and reports among government bodies, regulators, financial institutions and reporting entities. In this paper, we analyze the XBRL model from a perspective of data management to show why XBRL has created a new emerging vertical domain of database applications that offer opportunities and technical challenges for the database community. We outline the concept of an *XBRL repository* as a data hub to store and process collections of XBRL documents while maintaining document integrity based on XBRL semantics, providing document manageability, operational completeness and XBRL dictionary, query and business intelligence services. We then share our experiences of building the main components of the XBRL repository using the Oracle XML enabled RDBMS that leverages both XML and relational technologies as a backbone to host industrial strength XBRL applications. We also highlight technical challenges and potential solutions to each of the components. Finally we propose the potential of leveraging the XBRL data model concept for providing XML data normalization and XML having multi-hierarchy.

Visualizing Large-Scale RDF Data Using Subsets, Summaries, and Sampling in Oracle

Seema Sundara, Medha Atre, Vladimir Kolovski, Souripriya Das, Zhe Wu, Eugene Inseok Chong, Jagannathan Srinivasan; Oracle, USA

The paper addresses the problem of visualizing large scale RDF data via a *3-S* approach, namely, by using, 1) *Subsets*: to present only relevant data for visualisation; both static and dynamic subsets can be specified, 2) *Summaries*: to capture the essence of RDF data being viewed; summarized data can be expanded on demand thereby allowing users to create hybrid (summary-detail) fisheye views of RDF data, and 3) *Sampling*: to further optimize visualization of large-scale data where a representative sample suffices. The visualization scheme works with both asserted and inferred triples (generated using RDF(S) and OWL semantics). This scheme is implemented in Oracle by developing a plug-in for the Cytoscape graph visualization tool, which uses functions defined in a Oracle PL/SQL package, to provide fast and optimized access to Oracle Semantic Store containing RDF data. Interactive (Wikipedia 47 million triples and UniProt 700 million triples), and an OWL ontology (eClassOWl with a large class hierarchy including over 25,000 OWL classes, 5,000 properties, and 400,000 class-properties) demonstrates the effectiveness of our visualization scheme.

Industry Session 3: Query Optimization

Regency D, 08:30 – 10:00, Friday Chair: Vinayak Borkar

Incorporating Partitioning and Parallel Plans into the SCOPE Optimizer

Jingren Zhou, Per-Ake Larson, Ronnie Chaiken; Microsoft, USA

Massive data analysis on large clusters presents new opportunities and challenges for query optimization. Data partitioning is crucial to performance in this environment. However, data repartitioning is a very expensive operation so minimizing the number of such operations can yield very significant performance improvements. A query optimizer for this environment must therefore be able to reason about data partitioning including its interaction with sorting and grouping.

SCOPE is a SQL-like scripting language used at Microsoft for massive data analysis. A transformation-based optimizer is responsible for converting scripts into efficient execution plans for the Cosmos distributed computing platform. In this paper, we describe how reasoning about data partitioning is incorporated into the SCOPE optimizer. We show how

relational operators affect partitioning, sorting and grouping properties and describe how the optimizer reasons about and exploits such properties to avoid unnecessary operations. In most optimizers, consideration of parallel plans is an afterthought done in a postprocessing step. Reasoning about partitioning enables the SCOPE optimizer to fully integrate consideration of parallel, serial and mixed plans into the cost-based optimization. The benefits are illustrated by showing the variety of plans enabled by our approach.

Rule Profiling for Query Optimizers and Their Implications

Surajit Chaudhuri, Leo Giakoumakis, Vivek Narasayya, Ravishankar Ramamurthy; Microsoft, USA

Many modern optimizers use a transformation rule based framework. While there has been a lot of work on identifying new transformation rules, there has been little work focused on empirically evaluating the effectiveness of these transformation rules. In this paper we present the results of an empirical study of "profiling" transformation rules in Microsoft SQL Server using a diverse set of real world and benchmark query workloads. We also discuss the implications of these results for designing and testing query optimizers.

Data Desensitization of Customer Data for Use in Optimizer Performance Experiments

*Malu Castellanos*¹, *Bin Zhang*¹, *Ivo Jimenez*¹, *Perla Ruiz*², *Miguel Durazo*², *Umeshwar Dayal*¹, *Lily Jow*¹; ¹*HP*, USA; ²*University of Sonora, Mexico*

Improving the performance and functionality of database system optimizers requires experimentation on real customer data. Often these data are of sensitive nature and the only way to keep them is by applying a non-reversible transformation to obfuscate them. However, in order that the database optimizer generates exactly the same query plans as for the sensitive data, the transformation has to preserve the order and some important properties of the data distribution. Unfortunately, existing data obfuscation techniques do not preserve all of these properties and therefore are not applicable in this context. In this paper we present a Desensitizer tool that we have developed for optimizer performance experiments of HP's Neoview high availability data warehousing product. The tool is based on novel numeric and string desensitization algorithms which are agnostic to the database system. We explain the core concepts behind the algorithms, how they preserve the required data properties and important implementation considerations that were made. We present the architecture of the Desensitizer tool and results of the extensive validation that we conducted.

Demo Session 1A, 1B: Events, Streams, Services, Mashups and Search

Regency A, 10:30 - 12:00 Tuesday and 10:30 - 12:00 Wednesday

A Demonstration of the MaxStream Federated Stream Processing System

Irina Botan¹, Younggoo Cho², Roozbeh Derakhshan¹, Nihal Dindar¹, Ankush Gupta¹, Laura M. Haas³, Kihong Kim², Chulwon Lee², Girish Mundada⁴, Ming-Chien Shan⁴, Nesime Tatbul¹, Ying Yan⁵, Beomjin Yun², Jin Zhang⁵; ¹ETH Zürich, Switzerland; ²SAP, Korea; ³IBM, USA; ⁴SAP, USA; ⁵SAP, China

MaxStream is a federated stream processing system that seamlessly integrates multiple autonomous and heterogeneous Stream Processing Engines (SPEs) and databases. In this paper, we propose to demonstrate the key features of MaxStream using two application scenarios, namely the Sales Map & Spikes business monitoring scenario and the Linear Road Benchmark, each with a different set of requirements. More specifically, we will show how the MaxStream Federator can translate and forward the application queries to two different commercial SPEs (Coral8 and StreamBase), as well as how it does so under various persistency requirements.

E-Cube: Multi-Dimensional Event Sequence Processing Using Concept and Pattern Hierarchies

*Mo Liu*¹, Elke A. Rundensteiner¹, Kara Greenfield¹, Chetan Gupta², Song Wang², Ismail Ari³, Abhay Mehta²; ¹Worcester Polytechnic Institute, USA; ²HP, USA; ³Ozyegin University, Turkey

Many modern applications including tag based mass transit systems, RFID-based supply chain management systems and online financial feeds require special purpose event stream processing technology to analyze vast amounts of sequential multi-dimensional data available in real-time data feeds. Traditional online analytical processing (OLAP) systems are not designed for real-time pattern-based operations, while Complex Event Processing (CEP) systems are designed for sequence detection and do not support OLAP operations. We will demonstrate a novel E-Cube model that combines CEP and OLAP techniques for multidimensional event pattern analysis at different abstraction levels. A London transit scenario will be given to demonstrate the utility and performance of this proposed technology.

TargetSearch: A Ranking Friendly XML Keyword Search Engine

Ziyang Liu, Yichuan Cai, Yi Chen; Arizona State University, USA

This demo illustrates an XML search engine TargetSearch that addresses an open problem in XML keyword search: given relevant matches to keywords, how to compose query results properly so that they can be effectively ranked and easily digested by users. The approaches adopted in the literature generate either overwhelmingly large results or fragmentary results, both of which may cause the ranking schemes to be ineffective. Intuitively, each query has a search target and each result should contain exactly one instance of the search target along with its evidence. We developed TargetSearch which composes atomic and intact query results driven by users' search targets.

Efficient Fuzzy Type-Ahead Search in TASTIER

*Guoliang Li*¹, *Shengyue Ji*², *Chen Li*², *Jiannan Wang*¹, *Jianhua Feng*¹; ¹*Tsinghua University, China;* ²*University of California at Irvine, USA*

TASTIER is a research project on the new information-access paradigm called *type-ahead search*, in which systems find answers to a keyword query on-the-fly as users type in the query. In this paper we study how to support *fuzzy* type-ahead search in TASTIER. Supporting fuzzy search is important when users have limited knowledge about the exact representation of the entities they are looking for, such as people records in an online directory. We have developed and deployed several such systems, some of which have been used by many people on a daily basis. The systems received overwhelmingly positive feedbacks from users due to their friendly interfaces with the fuzzy-search feature. We describe the design and implementation of the systems, and demonstrate several such systems. We show that our efficient techniques can indeed allow this search paradigm to scale on large amounts of data.

MASS: A Multi-Facet Domain-Specific Influential Blogger Mining System

Yichuan Cai, Yi Chen; Arizona State University, USA

With rapid development of web 2.0 technology and e-business, bloggers play significant roles in the blogosphere as well as the external world. In particular, influential bloggers can bring great business values to modern enterprise. Despite that several systems for mining influential bloggers are available, they measure the influence of bloggers in general rather than domain specific, which is not applicable for real application requirements, such as business advertisement, personalized recommendation and so on. In this paper, we propose an effective model to mine the top-k influential bloggers according to their interest domains, the impact and attitude of the comments to their posts, as well as their authority in the network of page links. In this demonstration, we present MASS, an effective system

for mining influential bloggers. We will present the techniques in MASS, its experimental evaluation, as well as its applications.

Product EntityCube: A Recommendation and Navigation System for **Product Search**

*Jongwuk Lee*¹, *Seung-won Hwang*¹, *Zaiqing Nie*², *Ji-Rong Wen*²; ¹*POSTECH, Korea*; ²*Microsoft, China*

We demonstrate Product EntityCube, a product recommendation and navigation system. While the unprecedented scale of a product search portal enables to satisfy users with diverse needs, this scale also complicates product recommendation. Specifically, our target application poses a unique challenge of overcoming insufficient user profiles and feedbacks. To address this problem, we organize query results into clusters representing different user perceptions of similarity, and provide a navigational UI to handle personal interests. Specifically, we first discuss *hybrid object clustering* capturing diverse user interests from millions of Web pages and disambiguating different perceptions using feature-based similarity. We then discuss *skyline object ranking* to highlight interesting items at each cluster. Our demonstration illustrates how Product EntityCube can enrich user product shopping experiences.

Navigating Through Mashed-Up Applications with COMPASS

Daniel Deutch, Ohad Greenshpan, Tova Milo; Tel-Aviv University, Israel

Mashups integrate a set of complementary Web-services and data sources, often referred to as *mashlets*. We consider here a common scenario where the integrated mashlets are part of larger Web-applications, and their integration yields a set of inter-connected applications. We refer to them as *Mashed-up Applications* (abbr. *MashAPP*). The inter-connections between the mashlets enrich the individual Web-applications, but at the same time make the user navigation within them more intricate as actions in one application may affect others. To address this difficulty, we present *COMPASS*, a system that assists users in their navigation through *MashAPPs*. The system employs a novel top-k algorithm to propose users the most effective navigation paths for their specified goals. The suggestions are continually adapted to choices taken by the users while navigating.

GenerIE: Information Extraction Using Database Queries

Luis Tari¹, Phan Huy Tu¹, Jörg Hakenberg¹, Yi Chen¹, Tran Cao Son², Graciela Gonzalez¹, Chitta Baral¹; ¹Arizona State University, USA; ²New Mexico State University, USA

Information extraction systems are traditionally implemented as a pipeline of special-purpose processing modules. A major drawback of such an approach is that whenever a new extraction goal emerges or a module is improved, extraction has to be re-applied from scratch to the entire text corpus even though only a small part of the corpus might be affected. In this demonstration proposal, we describe a novel paradigm for information extraction: we store the parse trees output by text processing in a database, and then express extraction needs using queries, which can be evaluated and optimized by databases. Compared with the existing approaches, database queries for information extraction enable generic extraction and minimize reprocessing. However, such an approach also poses a lot of technical challenges, such as language design, optimization and automatic query generation. We will present the opportunities and challenges that we met when building GenerIE, a system that implements this paradigm.

Power-Aware Data Analysis in Sensor Networks

Daniel Klan¹, Katja Hose¹, Marcel Karnstedt², Kai-Uwe Sattler¹; ¹Ilmenau University of Technology, Germany; ²NUI Galway, Ireland

Sensor networks have evolved to a powerful infrastructure component for event monitoring in many application scenarios. In addition to simple filter and aggregation operations, an important task in processing sensor data is data mining — the identification of relevant information and patterns. Limited capabilities of sensor nodes in terms of storage and processing capacity, battery lifetime, and communication demand a power-efficient, preferably sensor-local processing. In this paper, we present AnduIN, a system for developing, deploying, and running in-network data mining tasks. The system consists of a data stream processing engine, a library of operators for sensor-local processing, a box-andarrow editor for specifying data mining tasks and deployment, a GUI providing the user with current information about the network and running queries, and an alerter notifying the user if a better query execution plan is available. At the demonstration site, we plan to show our system in action using burst detection as example application.

A View-Based Monitoring for Privacy-Aware Web Services

Hassina Meziane¹, Salima Benbernou¹, Aouda K. Zerdali¹, Mohand-Said Hacid², Mike Papazoglou³; ¹Université Paris Descartes, France; ²Université de Lyon, France; ³Tilburg University, The Netherlands

The demonstration addresses the problem of monitoring the compliance of privacy agreement that spells out a consumer's privacy rights and how their private information must be handled by the service provider. We present a Privacy Agreement Monitoring system, an easy-to-use, and an efficient tool for tightly controlling the private data usage flow dynamically in the area of web services. Some reasoning can be made upon the observations to enhance the compliance of the privacy agreement and enrich the knowledge on misuses.

Viewing a World of Annotations Through AnnoVIP

Konstantinos Karanasos, Spyros Zoupanos; INRIA, France

The proliferation of electronic content has notably lead to the apparition of large corpora of interrelated structured documents (such as HTML and XML Web pages) and semantic annotations (typically expressed in RDF), which further complement these documents. Documents and annotations may be authored independently by different users or programs. We present AnnoVIP, a peer-to-peer platform, capable of efficiently exploiting a multitude of annotated documents, based on innovative materialized views.

MashRank: Towards Uncertainty-Aware and Rank-Aware Mashups

Mohamed A. Soliman, Mina Saleeb, Ihab F. Ilyas; University of Waterloo, Canada

Mashups are situational applications that build data flows to link the contents of multiple Web sources. Often times, ranking the results of a mashup is handled in a materialize-then-sort fashion, since combining multiple data sources usually destroys their original rankings. Moreover, although uncertainty is ubiquitous on the Web, most mashup tools do not reason about or reflect such uncertainty.

We introduce MashRank, a mashup tool that treats ranking as a first-class citizen in mashup construction, and allows for rank-joining Web sources with uncertain information. To the best of our knowledge, no current tools allow for similar functionalities. MashRank encapsulates a new probabilistic model reflecting uncertainty in ranking, a set of techniques implemented as pipelined operators in mashup plans, and a probabilistic ranking infrastructure based on Monte-Carlo sampling.

T-Warehouse: Visual OLAP Analysis on Trajectory Data

Luca Leonardi¹, Gerasimos Marketos², Elias Frentzos², Nikos Giatrakos², Salvatore Orlando¹, Nikos Pelekis², Alessandra Raffaetà¹, Alessandro Roncato¹, Claudio Silvestri¹, Yannis Theodoridis²; ¹Università Ca' Foscari Venezia, Italy; ²University of Piraeus, Greece

Technological advances in sensing technologies and wireless telecommunication devices enable novel research fields related to the management of trajectory data. As it usually happens in the data management world, the challenge after storing the data is the implementation of appropriate analytics for extracting useful knowledge. However, traditional data warehousing systems and techniques were not designed for analyzing trajectory data. Thus, in this work, we demonstrate a framework that transforms the traditional data cube model into a trajectory warehouse. As a proof-of-concept, we implemented T-WAREHOUSE, a system that incorporates all the required steps for Visual Trajectory Data Warehousing, from trajectory reconstruction and ETL processing to Visual OLAP analysis on mobility data.

WikiAnalytics: Ad-Hoc Querying of Highly Heterogeneous Structured Data

Andrey Balmin¹, Emiran Curtmola²; ¹IBM, USA; ²University of California at San Diego, USA

Searching and extracting meaningful information out of highly heterogeneous datasets is a hot topic that received a lot of attention. However, the existing solutions are based on either rigid complex query languages (e.g., SQL, XQuery/XPath) which are hard to use without full schema knowledge, without an expert user, and which require up-front data integration. At the other extreme, existing solutions employ keyword search queries over relational databases [3], [1], [10], [9], [2], [11] as well as over semistructured data [6], [12], [17], [15] which are too imprecise to specify exactly the user's intent [16].

To address these limitations, we propose an alternative search paradigm in order to derive tables of precise and complete results from a very sparse set of heterogeneous records. Our approach allows users to disambiguate search results by navigation along conceptual dimensions that describe the records. Therefore, we cluster documents based on fields and values that contain the query keywords. We build a universal navigational lattice (UNL) over all such discovered clusters. Conceptually, the UNL encodes all possible ways to group the documents in the data corpus based on where the keywords hit.

We describe, WIKIANALYTICS, a system that facilitates data extraction from the Wikipedia infobox collection. WIKIANALYTICS provides a dynamic and intuitive interface that lets the average user explore the search results and construct homogeneous structured tables, which can be further queried and mashed up (e.g., filtered and aggregated) using the conventional tools.

SMARTINT: A System for Answering Queries Over Web Databases Using Attribute Dependencies

Ravi Gummadi, Anupam Khulbe, Aravind Kalavagattu, Sanil Salvi, Subbarao Kambhampati; Arizona State University, USA

Many web databases can be seen as providing partial and overlapping information about entities in the world. To answer queries effectively, we need to integrate the information about the individual entities that are fragmented over multiple sources. At first blush this is just the inverse of traditional database normalization problem — rather than go from a universal relation to normalized tables, we want to reconstruct the universal relation given the tables (sources). The standard way of reconstructing the entities will involve joining the tables. Unfortunately, because of the autonomous and decentralized way in which the sources are populated, they often do not have Primary Key – Foreign Key relations. While tables do share attributes, naive joins over these shared attributes can result in reconstruction of many spurious entities thus seriously compromising precision. Our system, SMARTINT is aimed at addressing the problem of data integration in such scenarios. Given a query, our system uses the Approximate Functional Dependencies(AFDs) to piece together a tree of relevant tables and schemas for joining them. The result tuples produced by our system are able to strike a favorable balance between precision and recall.

Demo Session 2A, 2B: Scalability, Design, Optimization and Visualization

Regency A, 13:30 - 15:00 Tuesday and 13:30 - 15:00 Wednesday

Mini-Me: A Min-Repro System for Database Software

Nicolas Bruno, Rimma V. Nehme; Microsoft, USA

Testing and debugging database software is often challenging and time consuming. A very arduous task for DB testers is finding a *min-repro* — the "simplest possible setup" that reproduces the original problem. Currently, a great deal of searching for min-repros is carried out manually using non-database-specific tools, which is both slow and error-prone. We propose to demonstrate a system, called *Mini-Me*, designed to ease and speed-up the task of finding min-repros in database-related products. *Mini-Me* employs several effective tools, including: the novel simplification transformations, the high-level language for creating search scripts and automation, the "record-and-replay" functionality, and the visualization of the search space and results. In addition to the standard *application mode*, the system can be interacted with in the *game mode*. The latter can provide an intrinsically motivating environment for developing successful search strategies by DB testers, which can be data-mined and recorded as patterns and used as recommendations for DB testers in the future. Potentially, a system like *Mini-Me* can save hours of time (for both customers and testers to isolate a problem), which could result in faster fixes and large cost savings to organizations.

I/O-Efficient Statistical Computing with RIOT

Yi Zhang, Weiping Zhang, Jun Yang; Duke University, USA

Statistical analysis of massive data is becoming indispensable to science, commerce, and society today. Such analysis requires efficient, flexible storage support and special optimization techniques. In this demo, we present RIOT (*R with I/O Transparency*), a system that extends R, a popular computing environment for statistical data analysis. RIOT makes R programs I/O-efficient in a way transparent to users. It features a flexible array storage manager and an optimization engine suitable for statistical and numerical operations. RIOT also seamlessly integrates with external database systems, offering additional opportunities for processing data that reside in databases by blurring the boundary between database and host-language processing. This demo will show how statistical computation can be effectively and efficiently handled by RIOT.

Interactive Physical Design Tuning

Nicolas Bruno, Surajit Chaudhuri; Microsoft, USA

In the last decade, automated physical design tuning became a relevant area of research. The process of tuning a workload became more flexible but also more complex, and getting the best design upfront became difficult. We propose a paradigm shift for physical design tuning, in which sessions are highly interactive, allowing DBAs to quickly try different options, identify problems, and obtain physical designs in an agile manner.

Visualizing Cost-Based XQuery Optimization

Andreas M. Weiner, Theo Härder, Renato Oliveira da Silva; University of Kaiserslautern, Germany

Developing a full-fledged cost-based XQuery optimizer is a fairly complex task. Nowadays, there is little knowledge concerning suitable cost formulae and optimization strategies for exploring and constraining the tremendously large search space. To allow for a fair assessment of different optimization strategies, physical algebra operators, and indexing approaches, we developed an extensible optimization framework. The framework is accompanied by a supportive visual explain tool that enables user interactions to refine the inspection and the comprehension of the query plans proposed. Using this tool, the optimizer can be dynamically reconfigured and the impact of different optimization strategies on the final query execution plan can be immediately visualized.

XML Reasoning Made Practical

Pierre Genevès¹, Nabil Layaïda²; ¹CNRS, France; ²INRIA, France

We present a tool for the static analysis of XPath queries and XML Schemas. The tool introduces techniques used in the field of verification (such as binary decision diagrams) in order to efficiently solve XPath query satisfiability, containment, and equivalence, in the presence of real-world XML Schemas. The tool can be used in query optimizers, in order to prove soundness of query rewriting. It can also be used in type-checkers and optimizing compilers that need to perform all kinds of compile-time analyses involving XPath queries and XML tree constraints.

TransScale: Scalability Transformations for Declarative Applications

Alexander Böhm, Erich Marth, Carl-Christian Kanne; University of Mannheim, Germany

The goal of the Demaq/TransScale system is to automate the distribution of complex application processes to large numbers of hosts. We implement distribution as a source-level transformation that turns the distribution-unaware application specification for a single host into a set of programs that can be executed on the various machines of a cluster.

Reverse Engineering Models from Databases to Bootstrap Application Development

Ankit Malpani¹, Philip A. Bernstein², Sergey Melnik³, James F. Terwilliger²; ¹IIT Madras, India; ²Microsoft, USA; ³Google, USA

Object-relational mapping systems have become often-used tools to provide application access to relational databases. In a database-first development scenario, the onus is on the developer to construct a meaningful object layer for the application because shipping tools, as ORM tools only ship database reverse-engineering tools that generate objects with a trivial one-to-one mapping. We built a tool, EdmGen++, that combines pattern-finding rules from conceptual modelling literature with configurable conditions that increase the likelihood that found patterns are semantically relevant. EdmGen++ produces a conceptual model with inheritance in Microsoft's Entity Data Model, which Microsoft's Entity Framework uses to support an executable object-to-relational mapping. The execution time of EdmGen++ on customer databases is reasonable for design-time.

HECATAEUS: Regulating Schema Evolution

*George Papastefanatos*¹, *Panos Vassiliadis*², *Alkis Simitsis*³, *Yannis Vassiliou*¹; ¹*National Technical University of Athens, Greece;* ²*University of Ioannina, Greece;* ³*HP, USA*

HECATAEUS is an open-source software tool for enabling impact prediction, what-if analysis, and regulation of relational database schema evolution. We follow a graph theoretic approach and represent database schemas and database constructs, like queries and views, as graphs. Our tool enables the user to create hypothetical evolution events and examine their impact over the overall graph before these are actually enforced on it. It also allows definition of rules for regulating the impact of evolution via (a) default values for all the nodes of the graph and (b) simple annotations for nodes deviating from the default behavior. Finally, HECATAEUS includes a metric suite for evaluating the impact of evolution events and detecting crucial and vulnerable parts of the system.

ROX: The Robustness of a Run-Time XQuery Optimizer Against Correlated Data

*Riham Abdel Kader*¹, *Peter A. Boncz*², *Stefan Manegold*², *Maurice van Keulen*¹; ¹*University of Twente, The Netherlands*; ²*CWI, The Netherlands*

We demonstrate ROX, a run-time optimizer of XQueries, that focuses on finding the best execution order of XPath steps and relational joins in an XQuery. The problem of join

ICDE 2010

ordering has been extensively researched, but the proposed techniques are still unsatisfying. These either rely on a cost model which might result in inaccurate estimations, or explore only a restrictive number of plans from the search space. ROX is developed to tackle these problems. ROX does not need any cost model, and defers query optimization to runtime intertwining optimization and execution steps. In every optimization step, sampling techniques are used to estimate the cardinality of unexecuted steps and joins to make a decision which sequence of operators to process next. Consequently, each execution step will provide updated and accurate knowledge about intermediate results, which will be used during the next optimization round. This demonstration will focus on: (*i*) illustrating the steps that ROX follows and the decisions it makes to choose a good join order, (*ii*) showing ROX's robustness in the face of data with different degree of correlation, (*iii*) comparing the performance of the plan chosen by ROX to different plans picked from the search space, (*iv*) proving that the run-time overhead needed by ROX is restricted to a small fraction of the execution time.

Symphony: A Platform for Search-Driven Applications

John C. Shafer, Rakesh Agrawal, Hady W. Lauw; Microsoft, USA

We present the design of Symphony, a platform that enables non-developers to build and deploy a new class of search-driven applications that combine their data and domain expertise with content from search engines and other web services. The Symphony prototype has been built on top of Microsoft's Bing infrastructure. While Symphony naturally makes use of the customization capabilities exposed by Bing, its distinguishing feature is the capability it provides to the application creator to combine their proprietary data and domain expertise with content obtained from Bing. They can also integrate specialized data obtained from web services to enhance the richness of their applications. Finally, Symphony is targeted at non-developers and provides cloud services for the creation and deployment of applications.

ProbClean: A Probabilistic Duplicate Detection System

George Beskales, Mohamed A. Soliman, Ihab F. Ilyas, Shai Ben-David, Yubin Kim; University of Waterloo, Canada

One of the most prominent data quality problems is the existence of duplicate records. Current data cleaning systems usually produce one clean instance (repair) of the input data, by carefully choosing the parameters of the duplicate detection algorithms. Finding the right parameter settings can be hard, and in many cases, perfect settings do not exist.

We propose ProbClean, a system that treats duplicate detection procedures as data processing tasks with uncertain outcomes. We use a novel uncertainty model that compactly encodes the space of possible repairs corresponding to different parameter settings. Prob-Clean efficiently supports relational queries and allows new types of queries against a set of possible repairs.

TransDec: A Spatiotemporal Query Processing Framework for Transportation Systems

Ugur Demiryurek, Farnoush Banaei-Kashani, Cyrus Shahabi; University of Southern California, USA

In this paper, we present *TransDec*, an end-to-end data-driven system which enables spatiotemporal queries in transportation systems with dynamic, real-time and historical data. TransDec fuses a variety of real-world spatiotemporal datasets including massive traffic sensor data, trajectory data, transportation network data, and point-of-interest data to create an immersive and realistic virtual model of a transportation system. With TransDec, we address the challenges in visualization, monitoring, querying and analysis of dynamic and large-scale transportation data in both time and space.

Provenance Browser: Displaying and Querying Scientific Workflow Provenance Graphs

Manish Kumar Anand¹, Shawn Bowers², Bertram Ludäscher¹; ¹University of California at Davis, USA; ²Gonzaga University, USA

This demonstration presents an interactive provenance browser for visualizing and querying data dependency (lineage) graphs produced by scientific workflow runs. The browser allows users to explore different views of provenance as well as to express complex and recursive graph queries through a high-level query language (QLP). Answers to QLP queries are lineage preserving in that queries return sets of lineage dependencies (denoting provenance graphs), which can be further queried and visually displayed (as graphs) in the browser. By combining provenance visualization, navigation, and query, the provenance browser can enable scientists to more easily access and explore scientific workflow provenance information.

Inconsistency Resolution in Online Databases

Yannis Katsis¹, Alin Deutsch¹, Yannis Papakonstantinou¹, Vasilis Vassalos²; ¹University of California at San Diego, USA; ²Athens University of Economics & Business, Greece

Shared online databases allow community members to collaboratively maintain knowledge. Collaborative editing though inevitably leads to inconsistencies as different members enter erroneous data or conflicting opinions. Ideally community members should be able to see and resolve these inconsistencies in a collaborative fashion. However most current online databases do not support inconsistency resolution. Instead they try to bypass the problem by either ignoring inconsistencies and treating data as if they were not conflicting or by requiring inconsistencies to be resolved outside the system.

To address this limitation, we propose Ricolla; an online database system that, by treating inconsistencies as first-class citizens, supports a natural workflow for the management of conflicting data. The system captures inconsistencies (so that community members can easily inspect them) and remains fully functional in their presence, thus enabling inconsistency resolution in an "as-you-go" fashion. Moreover it supports several schemes for the resolution of inconsistencies, allowing among others users to collaboratively resolve certain conflicts while disagreeing on others.

4th International Workshop on Ranking in Databases (DBRank'10)

Regency B, 08:30 - 14:30 Monday March 1st 2010

Program Co-Chairs:

- · Panagiotis G. Ipeirotis, New York University, USA
- · Amélie Marian, Rutgers University, USA

http://www.cs.uwaterloo.ca/conferences/dbrank/2010/

The Fourth International Workshop on Ranking in Databases (DBRank'10) 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 relational XML databases, text and multimedia databases, scientific information systems, biological databases, and so on. In particular, a large number of emerging applications require exploratory querying on such general-purpose or domain-specific databases: examples include users wishing to search bibliographic databases or catalogs of products such as homes, cars, cameras, restaurants, photographs, etc. To address the limitations of the traditional 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 multime-dia 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 user would be overwhelmed by too many unranked results. Furthermore, 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.

DBRank aims at providing more insight into supporting ranking in database systems and will be an interesting addition to ICDE 2010; the workshop will be a great venue for the many research groups working on ranking worldwide, with a unique opportunity to share their experience in supporting ranking in various database systems, from relational to semistructures and unstructured data; and on different levels from query formulation and preference modeling to query processing and optimization frameworks.

Keynote Speaker

Alon Halevy (Google)

Keynote Title: "Table Search"

Abstract: The Web contains hundreds of millions of high-quality structured data sets, in HTML tables, HTML lists and databases stored behind forms. Presenting these data sets in response to user queries is important for the quality of Web search. In addition, there are a growing number of contexts in which users would like to search specifically within a collection of structured data. Such search is hard because the typical signals used for text search do not apply as effectively to structured data. I will describe several scenarios in which table search is needed and contrast the treatments we've developed and continue to pursue at Google. In particular, I'll touch on our efforts to surface content from the Deep Web, our work on collecting and answering queries over a collection of 150 million HTML tables, and our work on supporting the creation and sharing of tables with Google Fusion Tables, a new service for collaborating on structured data on the Web

Speaker's Bio: Dr. Alon Halevy received his Bachelors degree in Computer Science and Mathematics from the Hebrew University in Jerusalem in 1988, and his Ph.D in Computer Science from Stanford University in 1993. From 1993 to 1997, Dr. Halevy was a principal member of technical staff at AT&T Bell Laboratories, and then at AT&T Laboratories. He joined the faculty of the Computer Science and Engineering Department at the University

of Washington in 1998. In 2004, Dr. Halevy founded Transformic Inc., a company that creates search engines for the deep web, content residing in databases behind web forms. Currently, he is at Google Inc. in Mountain View, California. Dr. Halevy was a Sloan Fellow (1999-2000), and received the Presidential Early Career Award for Scientists and Engineers (PECASE) in 2000. He serves on the editorial boards of the VLDB Journal, the Journal of Artificial Intelligence Research (currently, a member of the advisory committee), and ACM Transactions on Internet Technology. He served as the program chair for the ACM SIGMOD 2003 Conference, and has given several keynotes at top conferences.

Accepted Papers:

Subspace Similarity Search Using the Ideas of Ranking and Top-k Retrieval

Thomas Bernecker, Tobias Emrich, Franz Graf, Hans-Peter Kriegel, Peer Kröger, Matthias Renz, Erich Schubert, Arthur Zimek; LMU München, Germany

Efficient k-Nearest Neighbor Queries with the Signature Quadratic Form Distance

Christian Beecks, Merih Seran Uysal, Thomas Seidl; RWTH Aachen University, Germany

Top-k Pipe Join

Davide Martinenghi, Marco Tagliasacchi; Politecnico di Milano, Italy

On Novelty in Publish/Subscribe Delivery

Dimitris Souravlias, Marina Drosou, Kostas Stefanidis, Evaggelia Pitoura; University of Ioannina, Greece

Ranking for Data Repairs

Mohamed Yakout, Ahmed K. Elmagarmid, Jennifer Neville; Purdue University, USA

2nd IEEE Workshop on Information & Software as Services (WISS'10)

Regency D, March 1st 2010 08:30 - 18:00

Workshop General Chair: Divyakant Agrawal, *UC, Santa Barbara, USA* Workshop co-chairs:

- K. Selcuk Candan, Arizona State University, USA
- Wen-Syan Li, SAP Technology Lab, China

http://www.sap.com/china/about/ocs/en/wiss10/index.htm

Following the success of WISS'09 in Shanghai, China, the second IEEE workshop on "Information and Software as a Service" (WISS) will focus on the challenges associated with the design, deployment, and management of information and software as a service.

The high cost of creating and maintaining software and hardware infrastructures for delivering services to businesses led to a notable trend toward the use of third-party service providers, which rent out network presence, computation power, and data storage space to clients with infrastructural needs. These third party service providers can act as data stores as well as entire software suites for improved availability and system scalability, reducing small & medium businesses' burden of managing complex infrastructures. This is called information/application outsourcing or software as a service (SAAS). Emergence of enabling technologies, such as J2EE, .Net, XML, virtual machines, and web services contribute to this trend. Scientific Grid computing, on-line software services and business service networks are typical examples exploiting database and software as service paradigm.

KEYNOTE SPEAKER I

Mike Carey (UC Irvine)

Keynote Title: "Data Services: Past, Present, and Future"

Abstract: For many decades, data modeling and data access and update requirements played a central role in enterprise application development. Since the dawn of the SOA age in the early 2000's, however, data has gotten "lost in the suffle" and enterprises have struggled with the role of data in SOA. In response, IT vendors and analysts have contemplated - and in some cases, delivered - new software platforms for dealing specifically with "data services" or "information as a service". With the advent of cloud computing and the outsourcing of data management, the importance of defining and supporting data services is rapidly increasing.

This talk will take an opinionated look at the past, present, and future of data services. Different people mean different things by the term, so the talk will examine alternate definitions and use cases for data services including service-enablement of databases, data service composition from multiple data sources, and of course "cloud data management". The speaker will argue that things are currently getting messier by the minute in this area and will highlight several of the current key challenges and opportunities. (The talk will be based on the speaker's previous life at BEA Systems as well as recent brainstorming with Nicola Onose of UC Irvine and Michalis Petropoulos of SUNY Buffalo.)

KEYNOTE SPEAKER II

Jeff Hammerbacher (Cloudera, USA)

Keynote Title: "Open Questions for Building an Enterprise Data Platform on the Cloud"

Abstract: Innovations in hardware and software infrastructure are enabling new architectures for collecting, storing, and analyzing massive data sets. After a quick review of existing approaches, we'll propose some potential future scenarios and explore some of the research challenges engendered by the emergence of these new architectures.

Accepted Papers:

End-to-End Confidentiality for a Message Warehousing Service Using Identity-Based Encryption

Yuecel Karabulut¹, Harald Weppner¹, Ike Nassi¹, Anusha Nagarajan², Yash Shroff², Nishant Dubey², Tyelisa Shields²; ¹SAP, USA; ²Carnegie Mellon University, USA

The HiBench Benchmark Suite: Characterization of the MapReduce-Based Data Analysis

Shengsheng Huang, Jie Huang, Jinquan Dai, Tao Xie, Bo Huang; Intel, China

Towards Enterprise Software as a Service in the Cloud

*Jan Schaffner*¹, *Dean Jacobs*², *Benjamin Eckart*¹, *Jan Brunnert*¹, *Alexander Zeier*¹; ¹*HPI, Germany*; ²*SAP, Germany*

5th International Workshop on Self Managing Database Systems (SMDB'10)

Regency EF, March 1st 2010 08:30 - 18:00

Workshop Organizers:

- Shivnath Babu, Duke University, USA
- · Kai-Uwe Sattler, Technical University of Ilmenau, Germany

http://www.cs.duke.edu/smdb10/

Data management systems are growing rapidly in scale and complexity, while skilled administrators are becoming rarer and more expensive. Adding autonomic, or self-managing, capabilities to these systems promises easier use and maintenance. While considerable progress has been made in this direction, trends like cloud computing, virtualization, and software-as-a-service pose new challenges. Autonomic capabilities need to scale to hundreds of database nodes while taking economic factors into account.

The SMDB workshop brings together innovative researchers and practitioners to exchange ideas related to autonomic data management systems. SMDB'10 will be a one-day workshop where accepted papers are presented in an informal and interactive setting. Participation in the workshop is not limited to authors of accepted papers.

Previous workshops of the SMDB series focused on core topics in self-managing databases like physical design tuning, problem diagnosis and recovery, and database integration and protection. In addition to these core topics, the 2010 workshop seeks to broaden SMDB by soliciting submissions in emerging research areas like cloud computing, database testing, multitenant databases, large-scale storage systems, and datacenter administration.

Keynote Speaker:

Oliver Ratzesberger (Senior Director, Architecture and Operations, eBay)

Accepted Papers:

A Generic Auto-Provisioning Framework for Cloud Databases

*Jennie Rogers*¹, *Olga Papaemmanouil*², *Ugur Cetintemel*¹; ¹*Brown University, USA*; ²*Brandeis University, USA*

Adaptive Indexing for Relational Keys

Goetz Graefe, Harumi Kuno; HP, USA

Autonomic Workload Execution Control Using Throttling

*Wendy Powley*¹, *Patrick Martin*¹, *Mingyi Zhang*¹, *Paul Bird*², *Keith McDonald*²; ¹*Queen's University, Canada*; ²*IBM, Canada*

On the Use of Query-Driven XML Auto-Indexing

Karsten Schmidt, Theo Härder; University of Kaiserslautern, Germany

Statistics-Driven Workload Modeling for the Cloud

Archana Ganapathi, Yanpei Chen, Armando Fox, Randy Katz, David Patterson; University of California at Berkeley, USA

Vertical Partitioning of Relational OLTP Databases Using Integer Programming

Rasmus Resen Amossen; IT University of Copenhagen, Denmark

Automatic Tuning of the Multiprogramming Level in Sybase SQL Anywhere

*Mohammed Abouzour*¹, *Kenneth Salem*², *Peter Bumbulis*¹; ¹*Sybase iAnywhere, Canada;* ²*University of Waterloo, Canada*

Caching All Plans with Just One Optimizer Call

Debabrata Dash, Ioannis Alagiannis, Cristina Maier, Anastasia Ailamaki; EPFL, Switzerland

Towards Workload-Aware Self-Management: Predicting Significant Workload Shifts

Marc Holze, Ali Haschimi, Norbert Ritter; University of Hamburg, Germany

Panel: "Databases, MapReduce and the Cloud-Oh My! What's in it for the Administrator?"

Panelists:

- · Ashraf Aboulnaga (Univ. of Waterloo)
- Namit Jain (Facebook)
- Guy Lohman (IBM)
- Oliver Ratzesberger (eBay)
- · Benjamin Reed (Yahoo!)
- · Jingren Zhou (Microsoft)

2nd Workshop on Management and Mining of Uncertain Data (MOUND'10)

Regency C, March 1st 2010, 08:30-12:00

Program Chairs:

- · Heng Tao Shen, The University of Queensland, Australia
- · Graham Cormode, AT&T Labs-Research, USA

http://www.itee.uq.edu.au/~mound10/

Recently, uncertain data management and mining has become a critical issue in many real applications, such as sensor data monitoring, location-based services, object identification, and moving object search. Unlike exact data, uncertain data are often represented as a set of discrete samples or a probability density function, which presents new challenges for analyzing, querying, and mining the uncertain data effectively and efficiently. Following the success of the First International workshop on Management and mining Of UNcertain Data (MOUND'09), the Second MOUND'10 will continue to investigate key issues related to the data management and mining over uncertain data. Specifically, this forum welcomes contributions that explore uncertain data management issues such as data representation, various types of queries, and indexes. Additionally, the workshop hopes to attract work that studies the new data mining techniques of data cleaning, clustering, and classification over uncertain data, web and multimedia data.

INVITED TALK

Speaker: Dan Olteanu (University Lecturer, OUCL)

"A Toolbox of Query Evaluation Techniques for Probabilistic Databases"

In this talk I will discuss the problem of query evaluation in probabilistic databases and survey some of the existing techniques proposed by the database community. Although this problem is subsumed by general probabilistic inference, two fundamental aspects of databases, that is, (i) the separation of (very large) data and (small and fixed) query, and (ii) the use of mature relational query engines, can lead to more scalable techniques. I survey both exact and approximate query evaluation techniques. In case of exact evaluation, I discuss syntactical restrictions of the language of conjunctive queries with inequalities, under which the queries become tractable (in general, the problem is #P-hard). Relational query plans extended with efficient aggregation operators can be effectively used to evaluate such tractable queries. In case of intractable queries, exact techniques decompose the data-query instance into a tractable subinstance, which is solved as before, and a (usually much smaller) intractable subinstance that can be solved using AI inference techniques. Alternatively, intractable queries can be evaluated using (deterministic or randomized) approximation techniques with error guarantees.

Accepted Papers:

Constrained Frequent Itemset Mining from Uncertain Data Streams

Carson Kai-Sang Leung, Boyu Hao, Fan Jiang; University of Manitoba, Canada

Cleansing Uncertain Databases Leveraging Aggregate Constraints

*Haiquan Chen*¹, *Wei-Shinn Ku*¹, *Haixun Wang*²; ¹*Auburn University, USA*; ²*Microsoft, China*

U-DBSCAN : A Density-Based Clustering Algorithm for Uncertain Objects

Apinya Tepwankul, Songrit Maneewongwattana; King Mongkut's University of Technology Thonburi, Thailand

2nd International Workshop on New Trends in Information Integration (NTII'10)

March 5th 2010, 13:30 - 17:00 and March 6th 2010 08:30 - 15:15

Workshop Chair: Laura Haas, *IBM Almaden Research Center, USA* Program Committee Co-Chairs:

- · Zachary Ives, University of Pennsylvania, USA
- Manish A Bhide, *IBM Research, India*

http://www.cse.iitb.ac.in/~grajeev/ntii10/index.htm

Virtually every enterprise, scientific domain, or health care provider will assert that information integration is their most pressing information technology need. Despite the fact that research in data integration has been going on for over 20 years, we see few success stories from the real world. There are many reasons for this: perhaps predominantly that (1) integration encompasses a wide variety of tasks and domains, and there is a delicate balance between general solutions and domain-specific ones; and (2) general solutions typically require a combination of techniques from a range of communities, including databases, information retrieval, machine learning, and knowledge representation or Semantic Web. For instance, integrating contact center call transcripts with structured (transaction and profile) data in real-time requires efficient techniques which can work on noisy transcribed data, integrating Web data may need to deal with adversarial content providers, and integrating genetic data may require similarity matching on gene sequences. In recent years there has been a new emphasis on best-effort systems that combine automated approaches with user refinement or feedback, on integration techniques that combine the traditional stages of integration, and on using machine learning and other techniques with database concepts to address the needs of integration. These new approaches, generally targeting certain subclasses of the information integration problem, are highly promising.

The aim of the workshop is to encourage researchers from the information integration community to present novel issues and techniques related to applying information integration in different areas (especially in the context of integrating structured and unstructured data). The workshop will serve as a confluence of new ideas that will help drive research in the area of information integration from being 'generic' to being more focused, interactive, and realistic.

KEYNOTE SPEAKER 1

Dan Wolfson (IBM)

Keynote Title: "Business Information Management and Controls: Lessons from the Current Financial Crisis"

In the wake of the current financial crisis, the critical importance of Business Information Management and Controls has become increasingly evident. Many institutions around the world have been forced to evaluate their information systems for regulatory conformance, business efficiency and their ability to access risk/opportunity. As businesses continues to change their focus, merge together and divest unprofitable divisions, the challenges in providing information systems are many. The sheer complexity and scale is daunting. Political and economic realities must be accommodated. The regulatory requirements continue to evolve. To meet these challenges, a combination of software and engineering practices needs to be applied. In this talk we will explore some of the common information issues that have arisen and the technical and architectural approaches enlisted to improve the situation, We will focus on how to understand and record the meaning of information, its quality, and how it can be aggregated and shared across and organization. In addition to current practices, we will also highlight some of the key challenges and opportunities for the research community.

KEYNOTE SPEAKER 2

Craig Knoblock (USC)

Keynote Title: "Interactively Building Geospatial Mashups"

There are a number of tools and services available now for building mashups on the Web. However, many of the tools for constructing mashups reply on a widget paradigm, where users must select, customize, and connect widgets to build the desired application. While this approach does not require programming, the users must still understand programming concepts to successfully create a mashup. In this talk I describe our programming-bydemonstration approach to building mashups by example. Instead of requiring a user to select and customize a set of widgets, the user simply demonstrates the integration task by example. I will describe how this approach addresses the problems of extracting data from various sources, cleaning and modeling the extracted data, integrating the data across sources, and visualizing the integrated results in a geospatial context. We implemented these ideas in a system called Karma and evaluated Karma on a set of 20 users and showed that compared to other mashup construction tools, Karma allowed more of the users to successfully build mashups and made it possible to build these mashups significantly faster compared to using a widget-based approach. This research is joint work with Shubham Gupta, Pedro Szekely, and Rattapoom Tuchinda.

Accepted Papers:

A First Step Towards Integration Independence

*Laura M. Haas*¹, *Renée J. Miller*², *Donald Kossmann*³, *Martin Hentschel*³; ¹*IBM, USA*; ²*University of Toronto, Canada*; ³*ETH Zürich, Switzerland*

Towards Best-Effort Merge of Taxonomically Organized Data

David Thau¹, Shawn Bowers², Bertram Ludäscher¹; ¹University of California at Davis, USA; ²Gonzaga University, USA

Streaming Data Integration: Challenges and Opportunities

Nesime Tatbul; ETH Zürich, Switzerland

Partitioning Real-Time ETL Workflows

Alkis Simitsis, Chetan Gupta, Song Wang, Umeshwar Dayal; HP, USA

Midas for Government: Integration of Government Spending Data on Hadoop

Antonio Sala¹, Calvin Lin², Howard Ho³; ¹University of Modena & Reggio Emilia, Italy; ²University of California at Berkeley, USA; ³IBM, USA

Coordination of Data in Heterogenous Domains

Michael Lawrence, Rachel Pottinger, Sheryl Staub-French; University of British Columbia, Canada

BI-Style Relation Discovery Among Entities in Text

Wojciech M. Barczyński, Falk Brauer, Adrian Mocan, Marcus Schramm, Jan Froemberg; SAP, Germany

Profiling Linked Open Data with ProLOD

Christoph Böhm, Felix Naumann, Ziawasch Abedjan, Dandy Fenz, Toni Grütze, Daniel Hefenbrock, Matthias Pohl, David Sonnabend; HPI, Germany

Duplicate Detection in Probabilistic Data

*Fabian Panse*¹, *Maurice van Keulen*², *Ander de Keijzer*², *Norbert Ritter*¹; ¹*University of Hamburg, Germany;* ²*University of Twente, The Netherlands*

Complement Union for Data Integration

*Jens Bleiholder*¹, *Sascha Szott*², *Melanie Herschel*³, *Felix Naumann*¹; ¹*HPI, Germany*; ²*ZIB, Germany*; ³*Universität Tübingen, Germany*

1st International Workshop on Data Engineering Meets the Semantic Web (DESWEB'10)

Regency B, March 5th 2010 13:30 - 18:00 and March 6th 2010 08:30 - 12:00

General Chairs:

- Francesco Guerra, University of Modena and Reggio Emilia, Italy
- · Yannis Velegrakis, University of Trento, Italy

Panel Chair: Sonia Bergamaschi, University of Modena and Reggio Emilia, Italy

http://desweb2010.unimore.it/

The goal of the workshop is to bring together researchers and practitioners from the fields of Data Management and of the Semantic Web. It aims at investigating the new challenges that semantic web technologies have introduced and the ways through which these technologies can improve existing data management solutions.

KEYNOTE SPEAKER

Howard Ho (IBM Almaden)

Keynote Title: "Integrating Linked Open Data Sources about Financial Companies: Experience and Lessons Learned"

PANEL: TBA

Accepted Papers:

Semantic Flooding: Search Over Semantic Links

*Fausto Giunchiglia*¹, *Uladzimir Kharkevich*¹, *Alethia Hume*¹, *Piyatat Chatvorawit*²; ¹*University of Trento, Italy*; ²*Asian Institute of Technology, Thailand*

An Ontology-Based Retrieval System Using Semantic Indexing

*Soner Kara*¹, Özgür Alan¹, Orkunt Sabuncu¹, Samet Akpınar², Nihan K. Çiçekli², Ferda N. Alpaslan²; ¹Orbim Corp., Turkey; ²METU, Turkey

Keyword Based Search Over Semantic Data in Polynomial Time

*Paolo Cappellari*¹, *Roberto De Virgilio*², *Antonio Maccioni*², *Michele Miscione*²; ¹*University of Alberta, Canada;* ²*Università Roma Tre, Italy*

Towards Better Entity Resolution Techniques for Web Document Collections

Surender Reddy Yerva, Zoltán Miklós, Karl Aberer; EPFL, Switzerland

Processing Online News Streams for Large-Scale Semantic Analysis

*Miloš Krstajić*¹, *Florian Mansmann*¹, *Andreas Stoffel*¹, *Martin Atkinson*², *Daniel A. Keim*¹; ¹*University of Konstanz, Germany;* ²*EC Joint Research Centre, Italy*

DIVERSUM: Towards Diversified Summarisation of Entities in Knowledge Graphs

*Marcin Sydow*¹, *Mariusz Pikuła*¹, *Ralf Schenkel*²; ¹*Polish-Japanese Institute of Information Technology, Poland*; ²*Max-Planck Institute for Informatics, Germany*

The Entity Name System: Enabling the Web of Entities

Heiko Stoermer, Themis Palpanas, George Giannakopoulos; University of Trento, Italy

Ontology Alignment Argumentation with Mutual Dependency Between Arguments and Mappings

Paulo Maio, Nuno Silva; Politécnico do Porto, Portugal

Summarizing Ontology-Based Schemas in PDMS

*Carlos Eduardo Pires*¹, *Paulo Sousa*², *Zoubida Kedad*³, *Ana Carolina Salgado*²; ¹*UFCG, Brazil*; ²*UFPE, Brazil*; ³*UVSQ, France*
A Framework for Automatic Schema Mapping Verification Through Reasoning

*Paolo Cappellari*¹, *Denilson Barbosa*¹, *Paolo Atzeni*²; ¹*University of Alberta, Canada;* ²*Università Roma Tre, Italy*

Extensions to the Pig Data Processing Platform for Scalable RDF Data Processing Using Hadoop

Yusuke Tanimura, Akiyoshi Matono, Steven Lynden, Isao Kojima; AIST, Japan

Optimized Data Access for Efficient Execution of Semantic Services

Thorsten Möller, Heiko Schuldt; University of Basel, Switzerland

2nd International Workshop on Modeling, Managing and Mining of Evolving Social Networks (M3SN'10)

Regency C, March 6th 2010 08:30 - 12:00

Workshop Chairs:

- · Srikanta Bedathur, Max-Planck-Institut Informatik, Germany
- · Akshay Java, Microsoft, USA
- · Ralf Schenkel, Saarland University, Germany

http://www.mpi-inf.mpg.de/conferences/m3sn10/

Online social networking is quickly turning into an popular means of interacting with friends, sharing information, finding information as well as people with common interests, and, in general, a way to manage "personal spaces". Online Social networking services of all flavors have grown remarkably in a short span of time, with millions of users creating and sharing a vast amount of data ranging from blog entries, bookmarks, pictures to interactive games and personal interactions.

The amount of attention the research community has devoted to social networks has so far not kept up with their growth in popularity and overall importance. This needs to be addressed as social networks give rise to a number of new research challenges unique to them, such as modeling and exploiting their evolutionary dynamics, effective resource discovery within the variety of "social media" (photos, blogs, videos, maps, games, etc.) exploiting the number of interaction paths available, engineering of the mining algorithms required to deal with the underlying, heterogeneous data making up social networks, etc.

Adaptation of existing approaches in search and advertising to social networks is not straight-forward as well: the standard advertising models used in the context of sponsored search appear to break down when applied to social networks and searching for people, contacts or shared interests is very different from the search over documents studied in information retrieval. Exploiting the profusion of information within the social network requires deeper collaboration amongst research areas as diverse as graph theory to sociology to economics, with effective data engineering methods to assure scalability of the resulting methods.

In this workshop, we aim to address some of these open research challenges in modeling and mining of dynamic social networks.

Accepted Papers:

Privometer: Privacy Protection in Social Networks

Nilothpal Talukder, Mourad Ouzzani, Ahmed K. Elmagarmid, Hazem Elmeleegy, Mohamed Yakout; Purdue University, USA

On the Influence of Social Factors on Team Recommendations

Michele Brocco, Georg Groh, Christian Kern; Technische Universität München, Germany

Towards Discovery of Eras in Social Networks

*Michele Berlingerio*¹, *Michele Coscia*¹, *Fosca Giannotti*¹, *Anna Monreale*¹, *Dino Pedreschi*²; ¹*CNR-ISTI, Italy*; ²*University of Pisa, Italy*

Mining and Representing Recommendations in Actively Evolving Recommender Systems

Ira Assent; Aalborg University, Denmark

ICDE 2010 PhD Workshop

Beacon A, March 5th 2010 08:30 - 15:00

Program Co-Chairs:

- · Nikos Mamoulis, University of Hong Kong, China
- · Yannis Papakonstantinou, University of California, San Diego, USA
- · Timos Sellis, IMIS & NTU Athens, Greece

http://i.cs.hku.hk/~icde10ph/

The ICDE 2010 Ph.D. Workshop will bring together Ph.D. students working on topics related to the ICDE Conference series. The goal is to provide a forum for Ph.D. students to present ongoing research in a collaborative environment and share ideas with other young researchers in an international atmosphere. Participants will discuss their research ideas and results, and receive constructive feedback from an audience consisting of their peers as well as more senior researchers in their respective research areas. It will be an excellent opportunity for developing links and person-to-person networks which will be beneficial to our field and will be rewarding to the Ph.D. students. Furthermore, the workshop will include researchers, with experience in supervising and examining Ph.D. students, to provide feedback and advice to the participants.

Accepted Papers:

Maximizing Visibility of Objects

Muhammed Miah; University of Texas at Arlington, USA

Improving Product Search with Economic Theory

Beibei Li; New York University, USA

Toward Large Scale Data-Aware Search: Ranking, Indexing, Resolution and Beyond

Tao Cheng; University of Illinois at Urbana-Champaign, USA

Graphical Models for Dependencies and Queries in Uncertain Data

Ruiwen Chen; University of Ottawa, Canada

Privacy-Preserving Data Publishing

Ruilin Liu; Stevens Institute of Technology, USA

Flash-Enabled Database Storage

Ioannis Koltsidas; University of Edinburgh, UK A Database Server for Next-Generation Scientific Data Management Mohamed Y. Eltabakh; Purdue University, USA CareDB: A Context and Preference-Aware Location-Based Database System Justin J. Levandoski; University of Minnesota, USA **Graph Indexing for Reachability Queries** Hilmi Yıldırım; Rensselaer Polytechnic Institute, USA **Evaluating Path Queries Over Route Collections** Panagiotis Bouros; National Technical University of Athens, Greece Advances in Constrained Clustering ZiJie Qi; University of California at Davis, USA **Towards a Task-Based Search and Recommender Systems** Gabriele Tolomei: Università Ca' Foscari Venezia, Italy On Dynamic Data Clustering and Visualization Using Swarm Intelligence Esin Saka; University of Louisville, USA Fast Algorithms for Time Series Mining

Lei Li; Carnegie Mellon University, USA

ICDE 2010 COMMITTEES

ORGANIZING COMMITTEE

| General Chairs | Umeshwar Dayal, <i>HP Labs, USA</i> Vassilis Tsotras, <i>University of California at Riverside, USA</i> |
|---------------------------|---|
| Technical Program Chairs | Shahram Ghandeharizadeh, University of Southern California, USA Jayant R. Haritsa, Indian Institute of Science, India Gerhard Weikum, Max-Planck Institute for Informatics, Germany |
| Industrial Track Chairs | Mike Carey, University of California at Irvine, USA Fabio Casati, University of Trento, Italy Edward Chang, Google, China |
| Demo Chairs | Ioana Manolescu, INRIA, France Sharad Mehrotra, University of California at Irvine, USA |
| Panels Chairs | Anastasia Ailamaki, EPFL Lausanne, Switzerland Carlo Zaniolo, University of California at Los Angeles, USA |
| Tutorials/Seminars Chairs | Luis Gravano, <i>Columbia University, USA</i> Sunita Sarawagi, <i>IIT Bombay, India</i> |
| Workshops Chairs | Christian Jensen, Aalborg University, Denmark Renee Miller, University of Toronto, Canada |
| Proceedings Chairs | Feifei Li, Florida State University, USA Mirella Moro, Universidade Federal do Rio Grande do Sul, Brazil |
| Local Organization Chair | Chen Li, University of California at Irvine and BiMaple, USA |
| Finance Chair | Christof Bornhövd, SAP, USA |
| Publicity Chair | Roger Zimmermann, National University of Singapore, Singapore |

PROGRAM COMMITTEE

Technical Program Chairs

- · Shahram Ghandeharizadeh University of Southern California, USA
- · Jayant R. Haritsa Indian Institute of Science, India
- · Gerhard Weikum Max-Planck Institute for Informatics, Germany

Area PC Vice Chairs

- Social Networks and Personal Information *Sihem Amer-Yahia, Yahoo! Research*
- Semistructured Data Michael Benedikt, University of Oxford, UK
- Temporal, Spatial and Multimedia Data Michael H. Bohlen, Free University of Bozen-Bolzano, Italy
- Scientific Data -Susan B. Davidson, University of Pennsylvania, USA
- Text and Uncertain Data -Ihab Ilyas, University of Waterloo, Canada
- Indexing and Storage Masaru Kitsuregawa, University of Tokyo, Japan
- Distributed, Peer-to-Peer and Mobile Systems Qiong Luo, Hong Kong University of Science and Technology, China
- Systems, Experiments, Applications *Volker Markl, IBM Research, USA*

- Data Integration and Interoperability Felix Naumann, Hasso-Plattner-Institut, Germany
- Business Intelligence and Services Elke Rundensteiner, Worcester Polytechnic Institute, USA
- Query Processing and Optimization *S. Sudarshan, IIT Bombay, India*
- Streams and Sensor Networks *Kian-Lee Tan, National University of Singapore, Singapore*
- Data Mining Wei Wang, University of North Carolina at Chapel Hill, USA
- Privacy and Security Marianne Winslett, University of Illinois at Urbana-Champaign, USA
 Web Applications –
- Jeffrey Xu Yu, Chinese University of Hong Kong, China

Program Committee Members

- · Divy Agrawal University of California at Santa Barbara, USA
- Sanjay Agrawal *Microsoft, USA*
- Henrique Andrade *IBM*, *USA*
- · Arvind Arasu Microsoft, USA
- Walid Aref Purdue University, USA
- · Shivnath Babu Duke University, USA
- · Magdalena Balazinska University of Washington, USA
- Wolf-Tilo Balke TU Braunschweig, Germany
- · Andrey Balmin IBM, USA
- · Farnoush Banaei-Kashani University of Southern California, USA
- · Roberto Bayardo Google, USA
- Elisa Bertino Purdue University, USA
- Kevin Beyer IBM, USA
- · Jan-Geert Bex Hasselt University/Transnational University of Limburg, Belgium
- Peter Boncz CWI, The Netherlands
- · Angela Bonifati ICAR-CNR, Italy
- · Athman Bouguettaya CSIRO ICT Centre, Australia
- Daniele Braga Politecnico Milano, Italy
- Nico Bruno Microsoft, USA
- · Selcuk Candan Arizona State University, USA
- · Malu Castellanos HP, USA
- Tiziana Catarci University of Rome, Italy
- · Kaushik Chakrabarti Microsoft, USA
- · Chee Yong Chan National University of Singapore, Singapore
- · Kevin Chang University of Illinois at Urbana-Champaign, USA
- · Adriane Chapman MITRE, USA
- · Surajit Chaudhuri Microsoft, USA
- Lei Chen Hong Kong University of Science and Technology, China
- · Zhiyuan Chen University of Maryland at Baltimore County, USA
- · Reynold Cheng University of Hong Kong, China
- · Ravishankar Chinya University of California at Riverside, USA
- · Rada Chirkova North Carolina State University, USA
- · Laura Chiticariu IBM, USA
- Kajal Claypool MIT Lincoln Lab, USA
- · Chris Clifton Purdue University, USA
- · Sara Cohen Hebrew University, Israel
- · Sarah Cohen-Boulakia LRI, France
- Brian Cooper *Yahoo!, USA*
- Isabel Cruz University of Illinois at Chicago, USA

- · Nilesh Dalvi Yahoo!, USA
- · Gautam Das University of Texas Arlington, USA
- · Alex Delis University of Athens, Greece
- · Amol Deshpande University of Maryland, USA
- · Stefan Dessloch University of Kaiserslautern, Germany
- Alin Deutsch University of California at San Diego, USA
- · Yanlei Diao University of Massachusetts at Amherst, USA
- · Jens Dittrich University of Saarland, Germany
- Xin (Luna) Dong AT&T, USA
- Wenliang Du Syracuse University, USA
- Ahmed Elmagarmid Purdue University, USA
- Suzanne Embury University of Manchester, USA
- Vuk Ercegovac IBM, USA
- · Alexandre Evfimievski IBM, USA
- Wei Fan IBM, USA
- Wenfei Fan University of Edinburgh & Bell Labs, UK
- · Elena Ferrari University of Insubria, Italy
- · Juliana Freire University of Utah, USA
- · Keith Frikken Miami University Ohio, USA
- · Ada Waichee Fu Chinese University of Hong Kong, China
- · Avigdor Gal Technion, Israel
- · Cesar Galindo-Legaria Microsoft, USA
- · Johann Gamper University of Bolzano, Italy
- Sumit Ganguly IIT Kanpur, India
- · Minos Garofalakis Technical University of Crete, Greece
- · Lukasz Golab AT&T, USA
- · Goetz Graefe HP, USA
- · Le Gruenwald University of Oklahoma, USA
- · Torsten Grust University of Tübingen, Germany
- · Dimitrios Gunopulos University of California at Riverside, USA
- · Amarnath Gupta University of California at San Diego, USA
- · Ralf Hartmut Guting University of Hagen, Germany
- Marios Hadjieleftheriou AT&T, USA
- · Abdelkader Hameurlain IRIT, France
- · Jiawei Han University of Illinois at Urbana-Champaign, USA
- Wook-Shin Han Kyungpook National University, Korea
- · Vagelis Hristidis Florida International University, USA
- · Kien Hua University of Central Florida, USA
- Jun Huan University of Kansas, USA
- · Xiong Hui Rutgers University, USA
- · Yannis Ioannidis University of Athens, Greece
- · Panagiotis Ipeirotis NYU, USA
- · Yoshiharu Ishikawa Nagoya University, Japan
- · Dean Jacobs SAP, Germany
- H.V. Jagadish University of Michigan, USA
- · Christopher Jermaine University of Florida, USA
- Theodore Johnson AT&T, USA
- · Panos Kalnis KAUST, Saudi Arabia
- · Jaewoo Kang Korea University, Korea
- · Carl-Christian Kanne University of Mannheim, Germany
- Murat Kantarcioglu University of Texas Dallas, USA
- Norio Katayama NII, Japan
- · Raghav Kaushik Microsoft, USA
- · Daniel Keim University of Konstanz, Germany
- · Eamonn Keogh University of California at Riverside, USA
- George Kollios Boston University, USA

- Hank Korth Lehigh University, USA
- Nick Koudas University of Toronto, Canada
- Wei-Shinn Ku Auburn University, USA
- Harumi Kuno HP, USA
- · Alexandros Labrinidis University of Pittsburgh, USA
- · Laks Lakshmanan University of British Columbia, Canada
- Paul Larson Microsoft, USA
- Byung S. Lee University of Vermont, USA
- · Dik Lu n Lee Hong Kong University of Science and Technology, China
- · Dongwon Lee Penn State University, USA
- · Janice Lee National University of Singapore, Singapore
- Wang-Chien Lee Penn State University, USA
- · Kirsten LeFevre University of Michigan, USA
- · Toby Lehman IBM, USA
- Ulf Leser Humboldt University, Germany
- · Chen Li University of California at Irvine and BiMaple, USA
- · Jianzhong Li Harbin Institute of Technology, China
- Tao Li Florida International University, USA
- Wen-Syan Li SAP, China
- · Ee-Peng Lim Singapore Management University, Singapore
- Xuemin Lin University of New South Wales, Australia
- Bing Liu University of Illinois at Chicago, USA
- Jinze Liu University of Kentucky, USA
- Ling Liu Georgia Tech, USA
- · David Lomet Microsoft, USA
- Boon Thau Loo University of Pennsylvania, USA
- Yoelle Maarfek Google, Israel
- · Samuel Madden MIT, USA
- · Sanjay Madria University of Missouri-Rolla, USA
- · David Maier Portland State University, USA
- Nikos Mamoulis University of Hong Kong, China
- Stefan Manegold CWI, The Netherlands
- Amit Manjhi Google, USA
- · Ioana Manolescu INRIA, France
- Yun Mao AT&T, USA
- Wim Martens TU Dortmund, Germany
- · Claudia Bauzer Medeiros University of Campinas, Brazil
- · Sergey Melnik Google, USA
- Weiyi Meng Binghamton University, USA
- · Peter Mika Yahoo!, USA
- · Mukesh Mohania IBM, India
- Mohamed Mokbel University of Minnesota, USA
- John Mylopoulos University of Toronto, Canada
- Mario Nascimento University of Alberta, Canada
- Suman Nath Microsoft, USA
- · Jeffrey Naughton University of Wisconsin, USA
- · Shamkant Navathe Georgia Tech, USA
- · Thomas Neumann Max Planck Institute for Informatics, Germany
- · Alexandros Ntoulas Microsoft, USA
- · Frank Olken Lawrence Berkeley National Laboratory and NSF, USA
- Dan Olteanu University of Oxford, UK
- · Beng Chin Ooi National University of Singapore, Singapore
- · Ekow Otoo Lawrence Berkeley National Laboratory, USA
- · Sriram Padmanabhan IBM, USA
- · Dimitris Papadias Hong Kong University of Science and Technology, China
- Spiros Papadimitriou *IBM*, USA

- · Srinivasan Parthasarathy Ohio State University, USA
- · Jian Pei Simon Fraser University, Canada
- Antonella Poggi University of Rome, Italy
- · Neoklis Polyzotis University of California at Santa Cruz, USA
- Lucian Popa IBM, USA
- · Rachel Pottinger University of British Columbia, Canada
- Sunil Prabhakar Purdue University, USA
- · Weining Qian East China Normal University, China
- · Sriram Raghavan IBM, USA
- · Vijayshankar Raman IBM, USA
- · Berthold Reinwald IBM, USA
- · Mirek Riedewald Cornell University, USA
- · Doron Rotem LBNL Berkeley, USA
- Prasan Roy Aster Data Systems, USA
- · Pierangela Samarati University of Milan, Italy
- · Jorg Sander University of Alberta, Canada
- · Kai-Uwe Sattler TU Ilmenau, Germany
- · Ralf Schenkel University of Saarland, Germany
- · Peter Scheuermann Northwestern University, USA
- · Holger Schwarz University of Stuttgart, Germany
- · Bernhard Seeger University of Marburg, Germany
- Thomas Seidl RWTH Aachen, Germany
- · Timos Sellis IMIS and National Technical University of Athens, Greece
- · Pierre Senellart Telecom ParisTech, France
- · Cyrus Shahabi University of Southern California, USA
- · Jayavel Shanmugasundaram Yahoo!, USA
- · Mohamed Sharaf University of Toronto, Canada
- · Adam Silberstein Yahoo!, USA
- · Ambuj K. Singh University of California at Santa Barbara, USA
- Radu Sion Stony Brook University, USA
- · Yannis Sismanis IBM, USA
- · Richard Snodgrass University of Arizona, USA
- · Il-Yeol Song Drexel University, USA
- · Utkarsh Srivastava Yahoo!, USA
- · Robert Stevens University of Manchester, UK
- Dan Suciu University of Washington, USA
- · Wang-Chiew Tan University of California at Santa Cruz, USA
- · Xueyan Tang Nanyang Technological University, Singapore
- · Yufei Tao Chinese University of Hong Kong, China
- · Zahir Tari RMIT, Australia
- Nesime Tatbul ETH Zurich, Switzerland
- · Martin Theobald Max Planck Institute for Informatics, Germany
- Frank Tompa University of Waterloo, Canada
- Peter Triantafillou University of Patras, Greece
- Niki Trigoni University of Oxford, UK
- Juan Trujillo University of Alicante, Spain
- · Anthony Tung National University of Singapore, Singapore
- Deepak Turaga IBM, USA
- · Zografoula Vagena Microsoft, USA
- · Jaideep Vaidya Rutgers, USA
- Shivakumar Vaithyanathan IBM, USA
- · Maurice van Keulen University of Twente, The Netherlands
- · Vasilis Vassalos Athens University of Economics and Business, Greece
- Yannis Velegrakis University of Trento, Italy
- Stratis D. Viglas University of Edinburgh, UK
- Florian Waas Greenplum, USA

- · Haixun Wang IBM, USA
- Wei Wang University of New South Wales, Australia
- · Ji-Rong Wen Microsoft, China
- Kyu-Young Whang KAIST, Korea
- · Ouri Wolfson University of Illinois at Chicago, USA
- Kun-Lung Wu *IBM*, *USA*
- · Jianliang Xu Hong Kong Baptist University, China
- Xifeng Yan University of California at Santa Barbara, USA
- · Jun Yang Duke University, USA
- · Xiaoxin Yin Microsoft, USA
- Man Lung Yiu Hong Kong Polytechnic University, China
- · Haruo Yokota Tokyo Institute of Technology, Japan
- Philip Yu University of Illinois at Chicago, USA
- · Cong Yu Yahoo!, USA
- Ting Yu North Carolina State University, USA
- · Clement Yu University of Illinois at Chicago, USA
- · Mohammed Zaki Rensselaer Polytechnic Institute, USA
- · ChengXiang Zhai University of Illinois at Urbana-Champaign, USA
- · Baihua Zheng Singapore Management University, Singapore
- · Ding Zhou Facebook, USA
- · Jingren Zhou Microsoft, USA
- · Xiaofang Zhou University of Queensland, Australia
- · Yongluan Zhou University of Southern Denmark, Denmark

Industrial Program Committee

- · Dong Zhang Google, China
- Florian Daniel University of Trento, Italy
- · Halvard Skosgrud Thoughtworks, Australia
- · Honesty Young Intel, China
- · Kevin Wilkinson HP, USA
- · Zhen Hua Liu Oracle, USA
- Ming-Chien Shan SAP, China
- · Oliver Ratzesberger eBay, USA
- Paco Curbera IBM, USA
- Piero Fraternali *Politecnico di Milano, Italy*
- Stefan Tai University of Karlsruhe, Germany
- Till Westmann SAP, Germany
- · Todd Walter Teradata, USA
- · Yoshinori Hara Kyoto University, Japan

Demonstration Program Committee

- Nicolas Anciaux INRIA Rocquencourt, France
- Daniele Braga Politecnico di Milano, Italy
- · Kaushik Chakrabarti Microsoft, USA
- · Panos K. Chrysanthis University of Pittsburgh, USA
- · Sarah Cohen-Boulakia Université Paris XI, France
- · Elena Ferrari Università Insubria, Italy
- · Irini Fundulaki FORTH, Greece
- · Hakan Hacigumus NEC, USA
- Melanie Herschel University of Tübingen, Germany
- Ravi Jammalamadaka eBay Inc., USA
- · Dmitri Kalashnikov University of California at Irvine, USA
- · Christoph Koch Cornell University, USA

- Yiming Ma Nokia, USA
- · Benjamin Nguyen Université de Versailles Saint-Quentin, France
- Kjetil Norvag *NTNU, Norway*
- Michael Ortega Yahoo!, USA
- Nicoleta Preda Max Planck Institute for Informatics, Germany
- Praveen Rao University of Missouri, USA
- · Chinya Ravishankar University of California at Riverside, USA
- · Dawit Seid Teradata Corp., USA
- · Joshua Spiegel Oracle, USA
- Jianyong Wang Tsinghua University, China
- · Xiaofang Zhou University of Queensland, Australia

| | | | | | | | | | | | |
|------|--|------|------|--|------|------|------|--|--|------|------|
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |

| | | | | | | | | | | | |
|------|------|------|------|------|--|------|------|------|--|------|------|
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |

| | | | | | | | | | | | |
|------|--|------|------|--|------|------|------|--|--|------|------|
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |

| | | | | | | | | | | | |
|------|--|------|------|--|------|------|------|--|--|------|------|
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |