



Proceedings of the VLDB Endowment

Volume 4, No. 4 – January 2011

**Proceedings of the 37th International Conference on
Very Large Data Bases, Seattle, WA**

Editor-in-Chief:

H. V. Jagadish

Guest Editors:

José Blakeley, Joseph M. Hellerstein, Nick Koudas, Wolfgang Lehner, Sunita Sarawagi, Uwe Röhm

PVLDB – Proceedings of the VLDB Endowment

Volume 4, No. 4, January 2011.

The 37th International Conference on Very Large Data Bases, Seattle, WA.

Copyright 2011 VLDB Endowment

Permission to make digital or hard copies of portions of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for components of this work owned by others than VLDB Endowment must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists requires prior specific permission and/or a fee. Request permission to republish from PVLDB under email: info@vldb.org.

Volume 4, Number 4: VLDB 2011 Research Track Papers

Pages ii - vi and 208 - 266

ISSN 2150-8097, January 2011.

Additional copies only online at: portal.acm.org and www.vldb.org

TABLE OF CONTENTS

Front Matter

Copyright Notice	ii
Table of Contents	iii
PVLDB Review Board	iv

Letters

Letter from the Research Track Co-Chair.....	<i>Nick Koudas</i>	vi
--	--------------------	----

Research Papers

Large-Scale Collective Entity Matching	208
.....	<i>Vibhor Rastogi, Nilesh Dalvi, Minos Garofalakis</i>
Automatic Wrappers for Large Scale Web Extraction	219
.....	<i>Nilesh Dalvi, Ravi Kumar, Mohamed Soliman</i>
Fast Sparse Matrix-Vector Multiplication on GPUs: Implications for Graph Mining	231
.....	<i>Xintian Yang, Srinivasan Parthasarathy, P. Sadayappan</i>
Using Paxos to Build a Scalable, Consistent, and Highly Available Datastore	243
.....	<i>Jun Rao, Eugene J. Shekita, Sandeep Tata</i>
Fast Set Intersection in Memory	255
.....	<i>Bolin Ding and Arnd Christian König</i>

PVLDB REVIEW BOARD

VLDB 2011 General PC Co-Chairs

José Blakeley, Microsoft

Joe Hellerstein, University of California – Berkeley

VLDB 2011 Research Track Co-Chairs

Nick Koudas, University of Toronto and Sysomos Inc.

Wolfgang Lehner, Dresden University of Technology

Sunita Sarawagi, IIT Bombay

Reviewer

Ashraf Aboulnaga (University of Waterloo)

Sibel Adali (Rensselaer Polytechnic Institute)

Charu Aggarwal (IBM Watson Research Center)

Divyakant Agrawal (Univ. California, Santa Barbara)

Anastasia Ailamaki (EPFL Lausanne)

Gustavo Alonso (ETH Zurich)

Shivnath Babu (Duke University)

Roberto Bayardo (Google)

Elisa Bertino (Purdue University)

Peter Boncz (CWI, Netherlands)

Angela Bonifati (Icar-CNR)

Christof Bornhoevd (SAP Palo Alto)

Mike Cafarella (University of Washington)

K. Selcuk Candan (Arizona State University)

Malu Castellanos (HP Labs)

Tiziana Catarci (University of Rome)

Chee-Yong Chan (National University of Singapore)

Kevin Chang (University of Illinois, Urbana-Champaign)

Surajit Chaudhuri (Microsoft Research)

Rada Chirkova (North Carolina State University)

Jan Chomicki (University at Buffalo)

Chin-Wan Chung (Korea Advanced Institute of SaT)

Chris Clifton (Purdue University)

Christine Collet (Grenoble Institute of Technology)

Graham Cormode (AT&T Labs)

Gautam Das (University of Texas, Arlington)

Anish Das Sarma (Yahoo! Research)

Amol Deshpande (University of Maryland)

AnHai Doan (University of Wisconsin)

Xin Dong (AT&T Labs)

Alexandre Evfimievski (IBM Research)

Wenfei Fan (University of Edinburgh & Bell Labs)

Johann-Christoph Freytag (Humboldt-Universität Berlin)

Johannes Gehrke (Cornell University)

Rainer Gemulla (IBM Almaden Research Center)

Aristides Gionis (Yahoo! Research)

Goetz Graefe (HP Labs)

Torsten Grust (Universität Tübingen, Germany)

Giovanna Guerrini (University of Genova)

Dimitris Gunopulos (University of Athens, Greece)

Theo Haerder (University of Kaiserslautern)

Alon Halevy (Google)

Vagelis Hristidis (Florida International University)

Meichun Hsu (HP Labs, Palo Alto)

Ihab Ilyas (University of Waterloo)

Zachary Ives (University of Pennsylvania)

Dean Jacobs (SAP)

Christian Jensen (Aalborg University)

Chris Jermaine (University of Florida)

Raghav Kaushik (Microsoft Research)

Bettina Kemme (McGill University)
Eamonn Keogh (University of California, Riverside)
Martin Kersten (CWI)
Christoph Koch (Cornell University)
Flip Korn (AT&T Labs)
Donald Kossmann (ETH Zurich)
Alberto Laender (Federal University of Minas Gerais)
Dongwon Lee (Penn State University)
Kristen Lefevre (University of Michigan)
Chen Li (University of California, Irvine)
Bin Liu (University of Michigan)
David Lomet (Microsoft Research)
Samuel Madden (MIT)
Nikos Mamoulis (University of Hong Kong)
Ioana Manolescu (INRIA)
Claudia Medeiros (University of Campinas)
Sergey Melnik (Google)
Marco Mesiti (Universita degli Studi di Milano)
Chaitanya Mishra (Facebook Inc.)
Felix Naumann (University of Potsdam)
Raymond Ng (University of British Columbia)
Christopher Olston (Yahoo! Research)
Themis Palpanas (University of Trento)
Dimitris Papadias (Hong Kong University of SaT)
Stavros Papadopoulos (Chinese University of Hong Kong)
Stefano Paraboschi (University of Bergamo)
Jian Pei (Simon Fraser University)
Rachel Pottinger (University of British Columbia)
Vijayshankar Raman (IBM Almaden Research Centre)
Prakash Ramanan (Wichita State University)

PVLDB Information Director

Gerald Weber (University of Auckland)

Steering Committee

Serge Abiteboul, Peter Apers, Philip Bernstein, Elisa Bertino, Peter Buneman, Martin Kersten, Z. Meral Ozsoyuglu

Louisa Raschid (University of Maryland)
Kenneth Ross (Columbia University)
Elke Rundensteiner (Worcester Polytechnic Institute)
Yehoshua Sagiv (Hebrew University, Jerusalem)
Ken Salem (University of Waterloo)
Kai-Uwe Sattler (Ilmenau University of Technology)
Bernhard Seeger (University of Marburg)
Jayavel Shanmugasundaram (Yahoo! Research)
Kyuseok Shim (Seoul National University)
Divesh Srivastava (AT&T Labs)
Dan Suciu (University of Washington)
S. Sudarshan (IIT Bombay)
Kian-Lee Tan (National University of Singapore)
Val Tannen (University of Pennsylvania)
Jens Teubner (ETH Zurich)
Martin Theobald (Max-Planck-Institut für Informatik)
Frank Tompa (University of Waterloo)
Anthony Tung (National University of Singapore)
Patrick Valduriez (INRIA)
Wie Wang (University of North Carolina)
Gerhard Weikum (Max Planck Institute, Germany)
Yuqing Wu (Indiana University)
Fei Xu (Microsoft Search)
Sihem Yahia (Yahoo! Research)
Jun Yang (Duke University)
Cong Yu (Yahoo! Research)
Jefferey Yu (Chinese University of Hong Kong)
Ting Yu (North Carolina State University)
Xiaohui Yu (York University)
Justin Zobel (University of Melbourne)

VLDB 2011 Proceedings Chair

Uwe Röhm (University of Sydney)

LETTER FROM THE RESEARCH TRACK CO-CHAIR

I am happy to introduce the fourth issue of the PVLDB Journal Volume 4, consisting of papers accepted as part of the year-around, journal-style review process that will culminate with the presentation of these papers at the annual VLDB Conference to be held in Seattle, WA. This new publication process offers a faster review cycle and early dissemination of research results throughout the year.

This issue consists of five excellent papers reporting techniques for large scale entity matching, automated wrappers for web extraction, techniques for vector multiplication on GPUs, the design of highly available data stores and main memory set intersection techniques. We hope you find these papers thought provoking, inspiring and that you continue to consider PVLDB as the forum for publishing your best work.

Nick Koudas, University of Toronto and Sysomos Inc.

VLDB 2011 Research Track Co-Chair