



# Proceedings of the VLDB Endowment

Volume 14, No. 4 – December 2020

Editors in Chief:

**Xin Luna Dong and Felix Naumann**

Associate Editors:

**Alon Halevy, Anastasia Ailamaki, Angela Bonifati, Arun Kumar, Ashraf Aboulnaga,  
Eugene Wu, Floris Geerts, Graham Cormode, Jeffrey Xu Yu, Jiannan Wang, Jingren Zhou,  
Jorge Arnulfo Quiané Ruiz, Juliana Freire, Jun Yang, Martin Theobald, Nesime Tatbul,  
Paolo Papotti, Rainer Gemulla, Stefan Manegold, Stratos Idreos, Surajit Chaudhuri,  
Xuemin Lin, Yi Chen, Yufei Tao, Zachary Ives, Zhifeng Bao**

Publication Editors:

**Thorsten Papenbrock and Hannes Mühleisen**

PVLDB – Proceedings of the VLDB Endowment

Volume 14, No. 4, December 2020.

All papers published in this issue will be presented at the 47th International Conference on Very Large Data Bases, Copenhagen, Denmark, 2021.

## **Copyright 2020 VLDB Endowment**

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org).

Volume 14, Number 4, December 2020

Pages i – vii and 458 - 720

ISSN 2150-8097

Available at: <http://www.pvldb.org> and <https://dl.acm.org/journal/pvldb>

## TABLE OF CONTENTS

### Front Matter

Copyright Notice .....	i
Table of Contents .....	ii
PVLDB Organization and Review Board – Vol. 14 .....	iv
Editorial .....	vii

### Research Papers

Space- and Computationally-Efficient Set Reconciliation via Parity Bitmap Sketch (PBS) .....	458
<i>Long Gong, Ziheng Liu, Liang Liu, Jun Xu, Mitsunori Ogihara, Tong Yang</i>	
Astrid: Accurate Selectivity Estimation for String Predicates using Deep Learning .....	471
<i>Suraj Shetiya, Saravanan Thirumuruganathan, Nick Koudas, Gautam Das</i>	
Compact, Tamper-Resistant Archival of Fine-Grained Provenance .....	485
<i>Nan Zheng, Zack Ives</i>	
Rumble: Data Independence for Large Messy Data Sets .....	498
<i>Ingo Müller, Ghislain Fourny, Stefan Irimescu, Can Cikis, Gustavo Alonso</i>	
Capturing and querying fine-grained provenance of preprocessing pipelines in data science .....	507
<i>Adriane Chapman, Paolo Missier, Giulia Simonelli, Riccardo Torlone</i>	
Local Dampening: Differential Privacy for Non-numeric Queries via Local Sensitivity .....	521
<i>Victor Farias, Felipe Brito, Cheryl Flynn, Javam Machado, Subhabrata Majumdar, Divesh Srivastava</i>	
Mainlining Databases: Supporting Fast Transactional Workloads on Universal Columnar Data File Formats .....	534
<i>Tianyu Li, Matthew Butrovich, Amadou Ngom, Wan Shen Lim, Wes Mckinney, Andrew Pavlo</i>	
Accelerating Exact Constrained Shortest Paths on GPUs .....	547
<i>Shengliang Lu, Bingsheng He, Yuchen Li, Hao Fu</i>	
Towards an Efficient Weighted Random Walk Domination .....	560
<i>Songsong Mo, Zhifeng Bao, Ping Zhang, Zhiyong Peng</i>	
Scalable Mining of Maximal Quasi-Cliques: An Algorithm-System Codesign Approach .....	573
<i>Guimu Guo, Da Yan, Tamer Özsu, Zhe Jiang, Jalal Khalil</i>	
CALYPSO: Private Data Management for Decentralized Ledgers .....	586
<i>Eleftherios Kokoris Kogias, Enis Ceyhun Alp, Linus Gasser, Philipp Jovanovic, Ewa Syta, Bryan Ford</i>	
Stacked Filters: Learning to Filter by Structure.....	600
<i>Brian Hentschel, Stratos Idreos, Kyle Deeds</i>	
Maximizing Social Welfare in a Competitive Diffusion Model .....	613
<i>Prithu Banerjee, Laks V.s. Lakshmanan, Wei Chen</i>	

Understanding the Idiosyncrasies of Real Persistent Memory .....	626
<i>Shashank Gugnani, Arjun Kashyap, Xiaoyi Lu</i>	
Explaining Ranking Functions .....	640
<i>Amelie Marian, Abraham Gale</i>	
ConnectIt: A Framework for Static and Incremental Parallel Graph Connectivity Algorithms .....	653
<i>Laxman Dhulipala, Changwan Hong, Julian Shun</i>	
Quality of Sentiment Analysis Tools: The Reasons of Inconsistency .....	668
<i>Wissam Mammam Kouadri, Mourad Ouziri, Salima Benbernou, Karima Echihabi, Themis Palpanas, Iheb Benamor</i>	
Hindsight Logging for Model Training .....	682
<i>Rolando Garcia, Erick Liu, Vikram Sreekanti, Bobby Yan, Anusha Dandamudi, Joseph Gonzalez, Joseph Hellerstein, Koushik Sen</i>	
Scalable Structural Index Construction for JSON Analytics .....	694
<i>Lin Jiang, Junqiao Qiu, Zhijia Zhao</i>	
Efficient Join Algorithms For Large Database Tables in a Multi-GPU Environment .....	708
<i>Ran Rui, Hao Li, Yi-Cheng Tu</i>	

## **PVLDB ORGANIZATION AND REVIEW BOARD - Vol. 14**

### **Editors in Chief of PVLDB**

Xin Luna Dong (Amazon)  
Felix Naumann (HPI, University of Potsdam)

### **Associate Editors of PVLDB**

Ashraf Aboulnaga (Qatar Computing Research Institute,  
Hamad Bin Khalifa University)  
Anastasia Ailamaki (EPFL)  
Zhifeng Bao (RMIT University)  
Angela Bonifati (Lyon 1 University)  
Surajit Chaudhuri (Microsoft Research)  
Yi Chen (New Jersey Institute of Technology)  
Graham Cormode (University of Warwick)  
Juliana Freire (New York University)  
Floris Geerts (University of Antwerp)  
Rainer Gemulla (University of Mannheim)  
Alon Halevy (Facebook)  
Stratos Idreos (Harvard University)  
Zachary Ives (University of Pennsylvania)  
Arun Kumar (UC San Diego)  
Xuemin Lin (University of New South Wales)  
Stefan Manegold (CWI, Leiden University)  
Paolo Papotti (Eurecom)  
Jorge Arnulfo Quiané Ruiz (Technical University of Berlin)  
Yufei Tao (Chinese University of Hong Kong)  
Nesime Tatbul (Intel Labs and MIT)  
Martin Theobald (Université du Luxembourg)

Jiannan Wang (Simon Fraser University)  
Eugene Wu (Columbia University)  
Jun Yang (Duke University)  
Jeffrey Xu Yu (The Chinese University of Hong Kong)  
Jingren Zhou (Alibaba)

### **Publication Editors**

Thorsten Papenbrock (HPI, University of Potsdam)  
Hannes Mühleisen (CWI)

### **PVLDB Managing Editor**

Wolfgang Lehner (Dresden University of Technology)

### **PVLDB Advisory Committee**

Divesh Srivastava (AT&T Labs-Research)  
M. Tamer Özsu (University of Waterloo)  
Juliana Freire (New York University)  
Xin Luna Dong (Amazon)  
Peter Boncz (CWI)  
Lei Chen (Hong Kong University of Science and  
Technology)  
Graham Cormode (University of Warwick)  
Xiaofang Zhou (University of Queensland)  
Magdalena Balazinska (University of Washington)  
Fatma Ozcan (IBM Almaden)  
Felix Naumann (HPI, University of Potsdam)  
Peter Triantafillou (University of Warwick)

## Review Board

Abolfazl Asudeh (University of Illinois)  
Ahmed Eldawy (University of California, Riverside)  
Alan Fekete (University of Sydney)  
Alekh Jindal (Microsoft)  
Alex Ratner (University of Washington)  
Altigran da Silva (Universidade Federal do Amazonas)  
Anthony Tung (National University of Singapore)  
Antonios Deligiannakis (Technical University of Crete)  
Arijit Khan Nanyang (Technological University, Singapore)  
Arnau Prat (Sparsity Technologies)  
Ashwin Machanavajjhala (Duke University)  
Asterios Katsifodimos (Technical University of Delft)  
Avrilia Floratou (Microsoft)  
Babak Salimi (University of Washington)  
Badrish Chandramouli (Microsoft Research)  
Beng Chin Ooi (National University of Singapore)  
Bin Yang (Aalborg University)  
Boris Glavic (Illinois Institute of Technology)  
Byron Choi (Hong Kong Baptist University)  
Carlos Scheidegger (University of Arizona)  
Carsten Binnig (Technical University of Darmstadt)  
Ce Zhang (ETH Zurich)  
Chengfei Liu (Swinburne University of Technology)  
Chengkai Li (University of Texas at Arlington)  
Chris Jermaine (Rice University)  
Christian Bizer (University of Mannheim)  
Cong Yu (Google)  
Daisy Zhe Wang (University of Florida)  
Danica Porobic (Oracle)  
Davide Mottin (Aarhus University)  
Dimitris Papadias (Hong Kong University of Science and Technology)  
Dong Deng (Rutgers University)  
Eric Lo (Chinese University of Hong Kong)  
Essam Mansour (Concordia University)  
Fatma Ozcan (IBM Research)  
Flip Korn (Google)  
Florin Rusu (University of California, Merced)  
Fotis Psallidas (Microsoft)  
Francesco Bonchi (ISI Foundation)  
Gao Cong (Nanyang Technological University)  
George Fletcher (Technical University of Eindhoven)  
Georgia Koutrika (Athena Research Center)  
Hao Wei (Amazon)  
Heiko Mueller (New York University)  
Hong Cheng (Chinese University of Hong Kong)  
Hongzhi Wang (Harbin Institute of Technology)  
Hung Ngo (RelationalAI)  
Immanuel Trummer (Cornell University)  
Ingo Müller (ETH Zürich)  
Jana Giceva (Technical University of Munich)  
Jennie Rogers (Northwestern University)  
Jeong-Hyon Hwang (University at Albany, State University of New York)  
Jiaheng Lu (University of Helsinki)  
Jianliang Xu (Hong Kong Baptist University)

Jianxin Li (Deakin University)  
Jignesh Patel (University of Wisconsin)  
Johann Gamper (Free University of Bozen-Bolzano)  
Johannes Gehrke (Microsoft)  
Jonas Traub (Technical University of Berlin)  
Joy Arulraj (Georgia Tech)  
Ju Fan (Renmin University of China)  
K. Selçuk Candan (Arizona State University)  
Kai Zeng (Alibaba)  
Katja Hose (Aalborg University)  
Ken Salem (University of Waterloo)  
Kenneth A. Ross (Columbia University)  
Khuzaima Daudjee (University of Waterloo)  
Konstantinos Karanasos (Microsoft)  
Laurel Orr (Stanford University)  
Lei Chen (Hong Kong University of Science and Technology)  
Lei Zou (Peking University)  
Li Xiong (Emory University)  
Lu Chen (Aalborg University)  
Lu Qin (University of Technology Sydney)  
Manasi Vartak (Verta)  
Manos Athanassoulis (Boston University)  
Manos Karpathiotakis (Facebook)  
Marco Serafini (University of Massachusetts Amherst)  
Marcos Antonio Vaz Salles (University of Copenhagen)  
Mark Callaghan (MongoDB)  
Markus Weimer (Microsoft)  
Matei Zaharia (Stanford University, Databricks)  
Matteo Interlandi (Microsoft)  
Matthaios Olma (Microsoft Research)  
Meihui Zhang Beijing (Institute of Technology)  
Miao Qiao (University of Auckland)  
Michael H. Böhlen (University of Zurich)  
Michael Cafarella (University of Michigan)  
Mirek Riedewald (Northeastern University)  
Mohamed Mokbel (Qatar Computing Research Institute)  
Mohamed Sarwat (Arizona State University)  
Mohammad Sadoghi (University of California, Davis)  
Mourad Ouzzani (Qatar Computing Research Institute, Hamad Bin Khalifa University)  
Muhammad Aamir Cheema (Monash University)  
Murat Demirbas (University at Buffalo, SUNY)  
Nan Tang (Qatar Computing Research Institute, Hamad Bin Khalifa University)  
Nick Koudas (University of Toronto)  
Nikolaus Augsten (University of Salzburg)  
Norman May (SAP)  
Norman Paton (University of Manchester)  
Odysseas Papapetrou (Technical University of Eindhoven)  
Oliver A. Kennedy (University at Buffalo, SUNY)  
Paolo Merialdo (Roma Tre University)  
Paraschos Koutris (University of Wisconsin – Madison)  
Peter Boncz (Centrum Wiskunde & Informatica)  
Qin Zhang Indiana (University Bloomington)  
Raja Appuswamy (Eurecom)  
Ralf Schenkel (University of Trier)

Raul Castro Fernandez (University of Chicago)  
Raymond Chi-Wing Wong (Hong Kong University of Science and Technology)  
Reynold Cheng (The University of Hong Kong)  
Reza Akbarinia (INRIA)  
Ruoming Jin (Kent State University)  
Ryan Johnson (Amazon Web Services)  
S. Sudarshan (IIT Bombay)  
Sanjay Krishnan (University of Chicago)  
Saravanan Thirumuruganathan (Qatar Computing Research Institute, Hamad Bin Khalifa University)  
Sebastian Schelter (University of Amsterdam)  
Semih Salihoglu (University of Waterloo)  
Senjuti Basu Roy (New Jersey Institute of Technology)  
Shaoxu Song (Tsinghua University)  
Shimin Chen (Chinese Academy of Sciences)  
Sibo Wang (The Chinese University of Hong Kong)  
Silu Huang (Microsoft Research)  
Spyros Blanas (Ohio State University)  
Srikanth Kandula (Microsoft Research)  
Steffen Zeuch (German Research Centre for Artificial Intelligence - DFKI)  
Stijn Vansummeren (Université libre de Bruxelles)  
Sudeepa Roy (Duke University)  
Sudip Roy (Google)  
Tamer Özsu (University of Waterloo)  
Themis Palpanas (University of Paris, French University Institute - IUF)  
Tianzheng Wang (Simon Fraser University)  
Tingjian Ge (University of Massachusetts, Lowell)  
Thomas Heinis (Imperial College)  
Thomas Neumann (Technical University of Munich)  
Toon Calders (Universiteit Antwerpen)  
Umar Farooq Minhas (Microsoft Research)

Viktor Leis (Friedrich Schiller University Jena)  
Walid Aref (Purdue University)  
Wei-Shinn Ku (Auburn University)  
Weiren Yu (University of Warwick)  
Wendy Hui Wang (Stevens Institute of Technology)  
Wenjie Zhang (University of New South Wales)  
Wolfgang Gatterbauer (Northeastern University)  
Xi He (University of Waterloo)  
Xiang Zhao (National University of Defence Technology)  
Xiangyao Yu (University of Wisconsin – Madison)  
Xiaokui Xiao (National University of Singapore)  
Xiaolan Wang (Megagon Labs)  
Xin Cao (University of New South Wales)  
Xu Chu (Georgia Tech)  
Yannis Velegarakis (Utrecht University)  
Ye Yuan (Beijing Institute of Technology)  
Yeye He (Microsoft Research)  
Ying Zhang (University of Technology Sydney)  
Yinghui Wu (Case Western Reserve University)  
Yongjoo Park (University of Illinois at Urbana-Champaign)  
Yongxin Tong (Beihang University)  
Yu Yang (City University of Hong Kong)  
Yuchen Li (Singapore Management University)  
Yudian Zheng (Twitter)  
Yunjun Gao (Zhejiang University)  
Zechao Shang (University of Chicago)  
Zhenjie Zhang (Singapore R&D, Yitu Technology Ltd.)  
Zhewei Wei (Renmin University of China)  
Ziawasch Abedjan (Technical University of Berlin)  
Zoi Kaoudi (Technical University of Berlin)

## EDITORIAL

The fourth issue of PVLDB volume 14 embraces several fundamental and challenging topics in data management and its novel ramifications. In particular, it covers papers on graph processing and management, efficient and scalable analytical processes for various data formats, efficient join processing, provenance management, machine learning for data systems, among the other topics. This vast umbrella of topics reflects the novel directions of our community towards new areas (such as machine learning) as well as the importance of modern research flavors permeating a variety of data formats, query languages, and applications.

Reading the volume, we have been excited about the numerous papers on graph processing and social networks, which is one of the hottest topics nowadays. In particular, the issue faithfully reflects this trend and includes papers targeting advances on static and incremental parallel graph connectivity, as well as new algorithmic flavors and complexity bounds for quasi-clique mining, competitive influence maximization, efficient weighted random walk domination, and subgraph isomorphism for graph streams leveraging graph neural networks. All these papers show the inherent richness of our domain, which builds upon solid theoretical foundations coupled with timely applications and experimental assessment. Likewise, the topic of efficient and scalable data analytics is also majorly covered in this issue, with papers spanning from real-time analytics for time series data to highly efficient analytical processes and distributed query execution for JSON datasets. Moreover, join processing has been at the core of database management systems since its inception. This issue is not the exception with a paper showing the ability of our community to adapt to new applications and hardware requirements. As we are now living in an era where the norm is having multiple hardware devices and computing designs cohabiting in a single environment, it becomes crucial for our community to embrace this new technology in order to fully leverage its benefits. Last but not least, we are also excited to see a renaissance of data provenance under the new lens of data debugging, especially in data science pipelines. We observe an increasing interest in our community with two papers on this topic appearing in the current issue. These papers aim at understanding and auditing (data debugging) data science pipelines via data provenance. We expect bigger traction on data debugging in our community in the near future as the need for debugging (correcting, understanding, auditing, among others) data science pipelines is getting more prominent in current applications, especially in the life sciences, such as in the healthcare domain.

In these truly hard pandemic times, we really hope that this new issue will make your work time more enjoyable. Keep safe and healthy.

Angela Bonifati and Jorge-Arnulfo Quiané-Ruiz  
PVLDB Associate Editors