



# Proceedings of the VLDB Endowment

Volume 15, No. 10 – June 2022

Editors in Chief:  
**Fatma Özcan, Juliana Freire and Xuemin Lin**

Associate Editors:  
**Arun Kumar, Azza Abouzied, Beng Chin Ooi, Boris Glavic, Dan Suciu,  
Divyakant Agrawal, Eugene Wu, Georgia Koutrika, Ioana Manolescu,  
Jeffrey Xu Yu, Julia Stoyanovich, Jun Yang, K. Selçuk Candan,  
Khuzaima Daudjee, Laure Berti-Equille, Lei Chen, Mohamed Mokbel,  
Neoklis Polyzotis, Paolo Papotti, Peter Boncz, Sebastian Schelter,  
Sourav S Bhowmick, Surajit Chaudhuri, Themis Palpanas, Vanessa Braganholo,  
Viktor Leis, Wang-Chiew Tan, Wenjie Zhang, Wook-Shin Han, Xiaofang Zhou**

Publication Editors:  
**Lijun Chang and Xin Cao**

PVLDB – Proceedings of the VLDB Endowment

Volume 15, No. 10, June 2022.

All papers published in this issue will be presented at the 48th International Conference on Very Large Data Bases, Sydney, Australia, 2022.

## **Copyright 2022 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.

Volume 15, Number 10, June 2022

Pages i – vii and 1978 - 2296

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. 15 .....	iv

### Research Papers

VIP Hashing - Adapting to Skew in Popularity of Data on the Fly .....	1978
<i>Aarati Kakaraparthys, Jignesh Patel, Brian Kroth, Kwanghyun Park</i>	
Near-Data Processing in Database Systems on Native Computational Storage under HTAP Workloads .....	1991
<i>Tobias Vincon, Christian Knoedler, Leonardo Solis-vasquez, Arthur Bernhardt, Sajjad Tamimi, Lukas Weber, Florian Stock, Andreas Koch, Ilia Petrov</i>	
Hercules Against Data Series Similarity Search .....	2005
<i>Karima Echihabi, Panagiota Fatouropoulos, Kostas Zoumpatianos, Themis Palpanas, Houda Benbrahim</i>	
DISTILL: Low-Overhead Data-Driven Techniques for Filtering and Costing Indexes for Scalable Index Tuning .....	2019
<i>Tarique Siddiqui, Wentao Wu, Vivek Narasayya, Surajit Chaudhuri</i>	
Optimizing Machine Learning Inference Queries with Correlative Proxy Models.....	2032
<i>Zhihui Yang, Zuozhi Wang, Yicong Huang, Yao Lu, Chen Li, X. Sean Wang</i>	
Banyan: A Scoped Dataflow Engine for Graph Query Service.....	2045
<i>Li Su, Xiaoming Qin, Zichao Zhang, Rui Yang, Le Xu, Indranil Gupta, Wenyuan Yu, Zeng Kai, Jingren Zhou</i>	
Frequency Estimation Under Multiparty Differential Privacy: One-shot and Streaming .....	2058
<i>Ziyue Huang, Yuan Qiu, Ke Yi, Graham Cormode</i>	
Optimizing Inference Serving on Serverless Platforms .....	2071
<i>Ahsan Ali, Riccardo Pincioli, Feng Yan, Evgenia Smirni</i>	
Columnar Formats for Schemaless LSM-based Document Stores .....	2085
<i>Wail Y Alkowaileet, Michael Carey</i>	
Efficient Shortest Path Counting on Large Road Networks .....	2098
<i>Yu-xuan Qiu, Dong Wen, Lu Qin, Wentao Li, Ronghua Li, Ying Zhang</i>	
Towards Communication-efficient Vertical Federated Learning Training via Cache-enabled Local Update.....	2111
<i>Fangcheng Fu, Xupeng Miao, Jiawei Jiang, Huanran Xue, Bin Cui</i>	
DESIRE: An Efficient Dynamic Cluster-based Forest Indexing for Similarity Search in Multi-Metric Spaces.....	2121
<i>Yifan Zhu, Lu Chen, Yunjun Gao, Baihua Zheng, Pengfei Wang</i>	
ABC: Attributed Bipartite Co-clustering.....	2134
<i>Jung-hoon Kim, Kaiyu Feng, Gao Cong, Diwen Zhu, Wenyuan Yu, Chunyan Miao</i>	

Time Series Data Encoding for Efficient Storage: A Comparative Analysis in Apache IoTDB .....	2148
<i>Jinzhao Xiao, Yuxiang Huang, Changyu Hu, Shaoxu Song, Xiangdong Huang, Jianmin Wang</i>	
SA-LSM: Optimize Data Layout for LSM-tree Based Storage using Survival Analysis .....	2161
<i>Teng Zhang, Jian Tan, Xin Cai, Jianying Wang, Feifei Li, Jianling Sun</i>	
Improving Matrix-vector Multiplication via Lossless Grammar-Compressed Matrices .....	2175
<i>Paolo Ferragina, Giovanni Manzini, Travis Gagie, Dominik Köppl, Gonzalo Navarro, Manuel Striani, Francesco Tosoni</i>	
NFL: Robust Learned Index via Distribution Transformation .....	2188
<i>Shangyu Wu, Yufei Cui, Jinghuan Yu, Xuan Sun, Tei-wei Kuo, Chun Jason Xue</i>	
LEGOStore: A Linearizable Geo-Distributed Store Combining Replication and Erasure Coding .....	2201
<i>Hamidreza Zare, Viveck Cadambe, Bhuvan Urgaonkar, Nader Alfares, Praneet Soni, Chetan Sharma, Arif Merchant</i>	
Misinformation Mitigation under Differential Propagation Rates and Temporal Penalties .....	2216
<i>Michael Simpson, Laks V.s. Lakshmanan, Farnoosh Hashemi</i>	
Serving Deep Learning Models with Deduplication from Relational Databases.....	2230
<i>Lixi Zhou, Jiaqing Chen, Amitabh Das, Hong Min, Lei Yu, Ming Zhao, Jia Zou</i>	
Density-optimized Intersection-free Mapping and Matrix Multiplication for Join-Project Operations	2244
<i>Zichun Huang, Shimin Chen</i>	
Design Trade-offs for a Robust Dynamic Hybrid Hash Join .....	2257
<i>Shiva Jahangiri, Michael Carey, Johann-christoph Freytag</i>	
YeSQL: "You extend SQL" with Rich and Highly Performant User-Defined Functions in Relational Databases.....	2270
<i>Yannis E Foufoulas, Alkis Simitsis, Eleftherios Stamatogiannakis, Yannis Ioannidis</i>	
Magic Shapes for SHACL Validation.....	2284
<i>Shqiponja Ahmetaj, Bianca Löhnert, Magdalena Ortiz, Mantas Simkus</i>	

## **PVLDB ORGANIZATION AND REVIEW BOARD - Vol. 15**

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

Fatma Ozcan (Google)  
Juliana Freire (New York University)  
Xuemin Lin (University of New South Wales)

### **Associate Editors of PVLDB**

Arun Kumar (University of California, San Diego)  
Azza Abouzied (NYU Abu Dhabi)  
Beng Chin Ooi (NUS)  
Boris Glavic (Illinois Institute of Technology)  
Dan Suciu (University of Washington)  
Divyakant Agrawal (University of California, Santa Barbara)  
Eugene Wu (Columbia University)  
Georgia Koutrika (ATHENA)  
Ioana Manolescu (INRIA and Institut Polytechnique de Paris)  
Jeffrey Xu Yu (Chinese University of Hong Kong)  
Julia Stoyanovich (New York University)  
Jun Yang (Duke University)  
K. Seçuk Candan (Arizona State University)  
Khuzaima Daudjee (University of Waterloo)  
Laks Lakshmanan (The University of British Columbia)  
Laure Berti-Equille (IRD)  
Lei Chen (Hong Kong University of Science and Technology)  
Mohamed Mokbel (University of Minnesota, Twin Cities)  
Neoklis Polyzotis (Google)  
Paolo Papotti  
Peter Boncz (CWI)  
Sebastian Schelter (University of Amsterdam)  
Sharad Mehrotra (U.C. Irvine)  
Sourav S Bhowmick (Nanyang Technological University)

Surajit Chaudhuri (Microsoft Research)

Themis Palpanas (University of Paris)  
Vanessa Braganholo (Fluminense Federal University)  
Viktor Leis (Friedrich Schiller University Jena)  
Wang-Chiew Tan (Megagon Labs)  
Wenjie Zhang (University of New South Wales)  
Wook-Shin Han (POSTECH)  
Xiaofang Zhou (Hong Kong University of Science and Technology)

### **Publication Editors**

Lijun Chang (University of Sydney)  
Xin Cao (University of New South Wales)

### **PVLDB Managing Editor**

Wolfgang Lehner (Dresden University of Technology)

### **PVLDB Advisory Committee**

Felix Naumann (HPI)  
Juliana Freire (New York University)  
Xuemin Lin (U of New South Wales)  
Georgia Koutrika (Athena Research Center)  
Jun Yang (Duke University)  
Vanessa Braganholo (Universidade Federal Fluminense)  
Sourav S Bhowmick (Nanyang Technological University)  
Chris Jermaine (Rice University)  
Peter Triantafillou (University of Warwick)  
Xin Luna Dong (Facebook)  
Fatma Ozcan (Google)  
Lei Chen (Hong Kong University of S&T)  
Graham Cormode (University of Warwick)  
Divesh Srivastava (AT&T Labs-Research)  
Wolfgang Lehner (TU Dresden)

## Review Board

Abolfazl Asudeh (University of Michifan)  
Aécio Santos (New York University)  
Ahmed Eldawy (University of California, Riverside)  
Alexander Hall (RelationalAI)  
Alexander J Ratner (University of Washington)  
Aline Bessa (New York University)  
Alkis Simitsis (Athena Research Center)  
Altigran da Silva (Universidade Federal do Amazonas)  
AnHai Doan (University of Wisconsin-Madison)  
Anna Fariha (Microsoft)  
Anton Dignös (Free University of Bozen-Bolzano)  
Antonio Cavalcante Araujo Neto (University of Alberta)  
Arijit Khan (Nanyang Technological University)  
Arvind Arasu (Microsoft)  
Babak Salimi (University of California, San Diego)  
Bailu Ding (Microsoft Research)  
Bertram Ludaescher (University of Illinois)  
Bolong Zheng (Huazhong University of Science and Technology)  
Brandon Haynes (Gray Systems Lab, Microsoft)  
Byron Choi (Hong Kong Baptist University)  
Carlo Curino (Microsoft -- GSL)  
Carlos Scheidegger (The University of Arizona)  
Carsten Binnig (TU Darmstadt)  
Ce Zhang (ETH)  
Cheng Long (Nanyang Technological University)  
Chengfei Liu (Swinburne University of Technology)  
Chuan Lei (Instacart)  
Chunbin Lin (Amazon AWS)  
Curtis Dyreson (Utah State University)  
Dan Kifer (Pennsylvania State University)  
Dana M Van Aken (Carnegie Mellon University)  
Daniel Deutch (Tel Aviv University)  
Daniel Oliveira (UFF, Brazil)  
David Koop (Northern Illinois University)  
Davide Mottin (Aarhus University)  
Dong Xie (Penn State University)  
Eduardo Ogasawara (CEFET-RJ)  
Eleni Tzirita Zacharatou (TU Berlin)  
Fabio Porto (LNCC)  
Faisal Nawab (University of California at Irvine)  
Fan Zhang (Guangzhou University)  
Fateme Nargesian (University of Rochester)  
Fei Chiang (McMaster University)  
Florin Rusu (UC Merced)  
Floris Geerts (University of Antwerp)  
Fotis Psallidas (Microsoft)  
George Fletcher (Eindhoven University of Technology)  
George Papadakis (University of Athens)  
Gerhard Weikum (Max-Planck-Institut für Informatik)  
Germain Forestier (University of Haute Alsace)  
Guoliang Li (Tsinghua University)  
Haipeng Dai (Nanjing University)  
Harish Doraiswamy (Microsoft Research India)  
Heiko Mueller (DeepReason.ai)  
Herodotos Herodotou (Cyprus University of Technology)

Holger Pirk (Imperial College)  
Hongzhi Yin (The University of Queensland)  
Huiping Cao (New Mexico State University)  
Immanuel Trummer (Cornell)  
Ioana Manolescu (INRIA and Institut Polytechnique de Paris)  
Ippokratis Pandis (Amazon)  
Ishtiyaque Ahmad (University of California, Santa Barbara)  
Jae-Gil Lee (KAIST)  
Jana Giceva (TU Munich)  
Jeffrey Xu Yu (Chinese University of Hong Kong)  
Jens Teubner (TU Dortmund University)  
Jia Zou (Arizona State University)  
Jian Pei (Simon Fraser University)  
Jianguo Wang (Purdue University)  
Jiannan Wang (Simon Fraser University)  
Jianxin Li (Deakin University)  
Jianye Yang (Central South University)  
Jiwon Seo (Hanyang University)  
Johannes Gehrke (Microsoft)  
Jorge Arnulfo Quiane Ruiz (TU Berlin)  
Joseph Near (University of Vermont)  
Junhu Wang (Griffith University)  
Kaiping Zheng (National University of Singapore)  
Kangfei Zhao (The Chinese University of Hong Kong)  
Karima Echihabi (Mohammed VI Polytechnic University)  
Katja Hose (Aalborg University)  
Kenneth A Ross (Columbia University)  
Kostas Zoumpatianos (Snowflake Computing)  
Lei Zou (Peking University)  
Leopoldo Bertossi (Universidad Adolfo Ibanez)  
Li Xiong (Emory University)  
Lianke Qin (University of California, Santa Barbara)  
Lijun Chang (The University of Sydney)  
Lin Ma (Carnegie Mellon University)  
Long Yuan (Nanjing University of Science and Technology)  
Lu Qin (UTS)  
Luciano Barbosa (Universidade Federal de Pernambuco)  
Marcelo Arenas (Universidad Católica & IMFD)  
Maria Luisa Sapino (U. Torino)  
Matteo Lissandrini (Aalborg University)  
Matthias Boehm (Graz University of Technology)  
Matthias Renz (University of Kiel)  
Max Heimel (Snowflake)  
Maximilian Schleich (University of Washington)  
Meihui Zhang (Beijing Institute of Technology)  
Melanie Herschel (Universität Stuttgart)  
Michael Abebe (University of Waterloo)  
Min Xie (Instacart)  
Mirella M Moro (Universidade Federal de Minas Gerais)  
Mohamed Sarwat (Arizona State University)  
Mohammad Dashti (MongoDB)  
Mohammad Javad Amiri (University of Pennsylvania)  
Mohammad Sadoghi (University of California, Davis)  
Muhammad Aamir Cheema (Monash University)

Nikita Bhutani (Megagon Labs)  
Oliver A Kennedy (University at Buffalo, SUNY)  
Panos K. Chrysanthis (University of Pittsburgh)  
Paolo Missier (Newcastle University)  
Parth Nagarkar (NMSU)  
Paul Groth (University of Amsterdam)  
Peng CHENG (East China Normal University)  
Peter Pietzuch (Imperial College London)  
Pierangela Samarati (Universita delgi Studi di Milano)  
Pinar Karagoz (METU, Turkey)  
Pinar Tozun (IT University of Copenhagen)  
Prithu Banerjee (UBC)  
Raoni Lourenço (New York University)  
Raul Castro Fernandez (UChicago)  
Ravi Ramamurthy (Microsoft)  
Raymond Chi-Wing Wong (Hong Kong University of Science and Technology)  
Renata Borovica-Gajic (University of Melbourne)  
Reynold Cheng (The University of Hong Kong)  
Rui Mao (Shenzhen University)  
Ruoming Jin (Kent State University)  
Sai Wu (Zhejiang University)  
Sainyam Galhotra (University of Chicago)  
Sanjay Krishnan (University of Chicago)  
Sanjib Kumar Das (Google)  
Sayan Ranu (IIT Delhi)  
Sebastian Link (University of Auckland)  
Semih Salihoglu (University of Waterloo)  
Senjuti Basu Roy (New Jersey Institute of Technology)  
Sergey Melnik (Google)  
Shantanu Sharma (New Jersey Institute of Technology)  
Shaoxu Song (Tsinghua University)  
Sheng Wang (New York University)  
Shimin Chen (Chinese Academy of Sciences)  
Shumo Chu (University of California, Santa Barbara)  
Shweta Jain (University of Illinois, Urbana-Champaign)  
Sibo Wang (The Chinese University of Hong Kong)  
Srinivasan Keshav (University of Cambridge)  
Steffen Zeuch (DFKI GmbH)  
Steven E Whang (KAIST)  
Subarna Chatterjee (Harvard University)  
Sudip Roy (Google)  
Supun C Nakandala (University of California, San Diego)  
Tamer Özsu (University of Waterloo)  
Tarique A Siddiqui (Microsoft Research)  
Thomas Heinis (Imperial College)  
Thomas Neumann (TUM)  
Tianzheng Wang (Simon Fraser University)  
Tien Tuan Anh Dinh (Singapore University of Technology and Design)

Tilmann Rabl (HPI, University of Potsdam)  
Ting Yu (Qatar Computing Research Institute)  
Torben Bach Pedersen (Aalborg University)  
Torsten Grust (Universität Tübingen)  
Umar Farooq Minhas (Microsoft Research)  
Vasiliki Kalavri (Boston University)  
Verena Kantere (National Technical University of Athens)  
Victor Zakhary (Oracle)  
Vivek Narasayya (Microsoft Research)  
Vraj Shah (University of California, San Diego)  
Walid G Aref (Purdue)  
Wasay Abdul (Harvard)  
Wei Wang (Hong Kong University of Science and Technology (Guangzhou))  
Wei Lu (Renmin university of china)  
Weiren Yu (University of Warwick)  
Wen Hua (The University of Queensland)  
Wolfgang Lehner (TU Dresden)  
Xi He (University of Waterloo)  
Xiang Lian (Kent State University)  
Xiao Qin (IBM Research)  
Xiaofei Zhang (University of Memphis)  
Xiaokui Xiao (National University of Singapore)  
Xiaolan Wang (Megagon Labs)  
Xiaoyang Wang (Zhejiang Gongshang University)  
Xin Huang (Hong Kong Baptist University)  
Yael Amsterdamer (Bar-Ilan university)  
Yanyan Shen (Shanghai Jiao Tong University)  
Ye Yuan (Northeastern University)  
Yeye He (Microsoft Research)  
Yi Chen (NJIT)  
Yi Lu (MIT)  
Yikai Zhang (Chinese University of Hong Kong)  
Yinan Li (Microsoft Research)  
Ying Zhang (University of Technology Sydney)  
Yongxin Tong (Beihang University)  
Yuanyuan Zhu (Wuhan University)  
Yue Wang (Shenzhen Institute of Computing Sciences, Shenzhen University)  
Yufei Tao (Chinese University of Hong Kong)  
Yuliang Li (Megagon Labs)  
Yuncheng Wu (National University of Singapore)  
Yunjun Gao (Zhejiang University)  
Yuval Moskovitch (University of Michigan)  
Zhifeng Bao (RMIT University)  
Zhongle Xie (Zhejiang University)  
Zi Huang (University of Queensland)  
Ziawasch Abedjan (Leibniz Universität Hannover)  
Zohar Karnin (Amazon)  
Zsolt István (IT University of Copenhagen)

## **LETTER FROM THE EDITORS IN CHIEF**

We are pleased to present the tenth issue of PVLDB, Volume 15. This issue contains 24 papers in total including 21 regular research papers, 1 scalable data science (SDS) paper, and 2 experiments analysis & benchmark (EA&B) papers. A broad range of topics are covered in this issue including machine learning & applied AI for data management, graph data management, database performance, database engines, data privacy, and data quality.

For the first paper in this issue, Kakaraparthi et al. propose VIP hashing, a hash table method that uses lightweight mechanisms for learning the popularity distribution of keys and adapting to the skew in popularity on the fly. Next, Vincon et al. study near-data processing in database systems on native computational storage under HTAP Workloads. Echihabi et al. propose Hercules, a parallel tree-based technique for exact similarity search on massive disk-based data series collections. Siddiqui et al. present low-overhead data-driven techniques for filtering and costing indexes for scalable index tuning. Yang et al. optimize machine learning inference queries with correlative proxy models. Su et al. propose Banyan, a scoped dataflow engine for graph query service. Huang et al. study the problem of frequency estimation under both privacy and communication constraints, where the data is distributed among k parties. Ali et al. optimize inference serving on serverless platforms. Alkowaileet et al. propose techniques based on piggy-backing on Log-Structured Merge (LSM) tree events and tailored to document stores to store data in a columnar layout. Qiu et al. propose efficient solutions for shortest path counting on large road networks. Fu et al. introduce CELU-VFL, an efficient VFL training framework that exploits the local update technique to reduce the cross-party communication rounds. Zhu et al. propose DESIRE, an efficient dynamic cluster-based forest index for similarity search in multi-metric spaces. Kim et al. study the Attributed Bipartite Co-clustering (ABC) problem, which unifies the bipartite modularity optimization and attribute cohesiveness. Xiao et al. study the time series data encoding problem for efficient storage. Zhang et al. propose SA-LSM to use (S)urvival (A)nalysis for Log-Structure Merge Tree (LSM-tree) key-value (KV) stores. Ferragina et al. improve matrix-vector multiplication via lossless grammar-compressed matrices. Wu et al. propose NFL, a robust learned index via distribution transformation. Zare et al. design LEGOStore, an erasure coding (EC) based linearizable data store over geo-distributed public cloud data centers (DCs). Simpson et al. propose an information propagation model that captures important temporal aspects that have been well observed in the dynamics of fake news diffusion, in contrast with the diffusion of truth. Zhou et al. propose synergistic storage optimization techniques for duplication detection, page packing, and caching, to enhance database systems for model serving. Huang et al. present density-optimized intersection-free mapping and matrix multiplication for join-project operations. Jahangiri et al. design trade-offs for a robust dynamic hybrid hash join. Foufoula et al. present YeSQL, an SQL extension with rich UDF support along with a pluggable architecture to easily integrate it with either server-based or embedded database engines. Ahmetaj et al. present magic shapes for SHACL validation.

All the papers in this issue will be presented at the 48th International Conference on Very Large Data Bases, 2022, in Sydney. We sincerely thank all the authors for submitting their work and all the reviewers for their outstanding service in reviewing the submissions. We hope that the reader will find this volume enjoyable.

Fatma Özcan, Juliana Freire and Xuemin Lin  
Editors-in-Chief of PVLDB Volume 15  
Program Chairs for VLDB 2022