

# **12th USENIX Symposium on Operating Systems Design and Implementation (OSDI'16)**

Savannah, Georgia, USA  
2 - 4 November 2016

ISBN: 978-1-7138-0466-6

**Printed from e-media with permission by:**

Curran Associates, Inc.  
57 Morehouse Lane  
Red Hook, NY 12571



**Some format issues inherent in the e-media version may also appear in this print version.**

Copyright© (2016) by Usenix Association  
All rights reserved.

Printed with permission by Curran Associates, Inc. (2020)

For permission requests, please contact Usenix Association  
at the address below.

Usenix Association  
2560 Ninth Street, Suite 215  
Berkeley, California, 94710

<https://www.usenix.org/>

**Additional copies of this publication are available from:**

Curran Associates, Inc.  
57 Morehouse Lane  
Red Hook, NY 12571 USA  
Phone: 845-758-0400  
Fax: 845-758-2633  
Email: [curran@proceedings.com](mailto:curran@proceedings.com)  
Web: [www.proceedings.com](http://www.proceedings.com)

**OSDI '16:**  
**12th USENIX Symposium on Operating Systems**  
**Design and Implementation**  
**Savannah, GA, USA**

Message from the Program Co-Chairs. . . . . vii

**Wednesday, November 2, 2016**

**Operating Systems I**

**Push-Button Verification of File Systems via Crash Refinement. . . . .1**

Helgi Sigurbjarnarson, James Bornholt, Emina Torlak, and Xi Wang, *University of Washington*

**Intermittent Computation without Hardware Support or Programmer Intervention . . . . .17**

Joel Van Der Woude, *Sandia National Laboratories*; Matthew Hicks, *University of Michigan*

**Machine-Aware Atomic Broadcast Trees for Multicores . . . . .33**

Stefan Kaestle, Reto Achermann, Roni Haecki, Moritz Hoffmann, Sabela Ramos, and Timothy Roscoe, *ETH Zurich*

**Light-Weight Contexts: An OS Abstraction for Safety and Performance . . . . .49**

James Litton, *University of Maryland, College Park and Max Planck Institute for Software Systems (MPI-SWS)*; Anjo Vahldiek-Oberwagner, Eslam Elnikety, and Deepak Garg, *Max Planck Institute for Software Systems (MPI-SWS)*; Bobby Bhattacharjee, *University of Maryland, College Park*; Peter Druschel, *Max Planck Institute for Software Systems (MPI-SWS)*

**Cloud Systems I**

**Altruistic Scheduling in Multi-Resource Clusters. . . . .65**

Robert Grandl, *University of Wisconsin—Madison*; Mosharaf Chowdhury, *University of Michigan*; Aditya Akella, *University of Wisconsin—Madison*; Ganesh Ananthanarayanan, *Microsoft*

**GRAPHENE: Packing and Dependency-Aware Scheduling for Data-Parallel Clusters . . . . .81**

Robert Grandl, *Microsoft and University of Wisconsin—Madison*; Srikanth Kandula and Sriram Rao, *Microsoft*; Aditya Akella, *Microsoft and University of Wisconsin—Madison*; Janardhan Kulkarni, *Microsoft*

**Firmament: Fast, Centralized Cluster Scheduling at Scale . . . . .99**

Ionel Gog, *University of Cambridge*; Malte Schwarzkopf, *MIT CSAIL*; Adam Gleave and Robert N. M. Watson, *University of Cambridge*; Steven Hand, *Google, Inc.*

**Morpheus: Towards Automated SLOs for Enterprise Clusters. . . . .117**

Sangeetha Abdu Jyothi, *Microsoft and University of Illinois at Urbana—Champaign*; Carlo Curino, Ishai Menache, and Shravan Matthur Narayanamurthy, *Microsoft*; Alexey Tumanov, *Microsoft and Carnegie Mellon University*; Jonathan Yaniv, *Technion—Israel Institute of Technology*; Ruslan Mavlyutov, *Microsoft and University of Fribourg*; Íñigo Goiri, Subru Krishnan, Janardhan Kulkarni, and Sriram Rao, *Microsoft*

**Transactions and Storage**

**The SNOW Theorem and Latency-Optimal Read-Only Transactions. . . . .135**

Haonan Lu, *University of Southern California*; Christopher Hodsdon, *University of Southern California*; Khiem Ngo, *University of Southern California*; Shuai Mu, *New York University*; Wyatt Lloyd, *University of Southern California*

**Correlated Crash Vulnerabilities . . . . .151**

Ramnatthan Alagappan, Aishwarya Ganesan, Yuvraj Patel, Thanumalayan Sankaranarayanan Pillai, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau, *University of Wisconsin—Madison*

(Wednesday, November 2, continues on next page)

**Incremental Consistency Guarantees for Replicated Objects** .....169  
Rachid Guerraoui, Matej Pavlovic, and Dragos-Adrian Seredinschi, *École Polytechnique Fédérale de Lausanne (EPFL)*

**FaSST: Fast, Scalable and Simple Distributed Transactions with Two-Sided (RDMA) Datagram RPCs** ...185  
Anuj Kalia, *Carnegie Mellon University*; Michael Kaminsky, *Intel Labs*; David G. Andersen, *Carnegie Mellon University*

## Networking

**NetBricks: Taking the V out of NFV** .....203  
Aurojit Panda and Sangjin Han, *University of California, Berkeley*; Keon Jang, *Google*; Melvin Walls and Sylvia Ratnasamy, *University of California, Berkeley*; Scott Shenker, *University of California, Berkeley, and International Computer Science Institute*

**Efficient Network Reachability Analysis Using a Succinct Control Plane Representation** .....217  
Seyed K. Fayaz and Tushar Sharma, *Carnegie Mellon University*; Ari Fogel, *Intentionet*; Ratul Mahajan, *Microsoft Research*; Todd Millstein, *University of California, Los Angeles*; Vyas Sekar, *Carnegie Mellon University*; George Varghese, *University of California, Los Angeles*

**Simplifying Datacenter Network Debugging with PathDump** .....233  
Praveen Tammana, *University of Edinburgh*; Rachit Agarwal, *Cornell University*; Myungjin Lee, *University of Edinburgh*

**Network Requirements for Resource Disaggregation** .....249  
Peter X. Gao, Akshay Narayan, Sagar Karandikar, Joao Carreira, and Sangjin Han, *University of California, Berkeley*; Rachit Agarwal, *Cornell University*; Sylvia Ratnasamy, *University of California, Berkeley*; Scott Shenker, *University of California, Berkeley, and International Computer Science Institute*

## Thursday, November 3, 2016

### Graph Processing and Machine Learning

**TensorFlow: A System for Large-Scale Machine Learning** .....265  
Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng, *Google Brain*

**Exploring the Hidden Dimension in Graph Processing** .....285  
Mingxing Zhang, Yongwei Wu, and Kang Chen, *Tsinghua University*; Xuehai Qian, *University of Southern California*; Xue Li and Weimin Zheng, *Tsinghua University*

**Gemini: A Computation-Centric Distributed Graph Processing System** .....301  
Xiaowei Zhu, Wenguang Chen, and Weimin Zheng, *Tsinghua University*; Xiaosong Ma, *Hamad Bin Khalifa University*

**Fast and Concurrent RDF Queries with RDMA-Based Distributed Graph Exploration** .....317  
Jiaxin Shi, Youyang Yao, Rong Chen, and Haibo Chen, *Shanghai Jiao Tong University*; Feifei Li, *University of Utah*

### Languages and Software Engineering

**RE<sup>X</sup>: A Development Platform and Online Learning Approach for Runtime Emergent Software Systems** ...333  
Barry Porter, Matthew Grieves, Roberto Rodrigues Filho, and David Leslie, *Lancaster University*

**Yak: A High-Performance Big-Data-Friendly Garbage Collector** .....349  
Khanh Nguyen, Lu Fang, Guoqing Xu, and Brian Demsky, *University of California, Irvine*; Shan Lu, *University of Chicago*; Sanazsadat Alamian, *University of California, Irvine*; Onur Mutlu, *ETH Zurich*

**Shuffler: Fast and Deployable Continuous Code Re-Randomization** .....367  
David Williams-King and Graham Gobieski, *Columbia University*; Kent Williams-King, *University of British Columbia*; James P. Blake and Xinhao Yuan, *Columbia University*; Patrick Colp, *University of British Columbia*; Michelle Zheng, *Columbia University*; Vasileios P. Kemerlis, *Brown University*; Junfeng Yang, *Columbia University*; William Aiello, *University of British Columbia*

**Don't Get Caught in the Cold, Warm-up Your JVM: Understand and Eliminate JVM Warm-up Overhead in Data-Parallel Systems** .....383  
David Lion and Adrian Chiu, *University of Toronto*; Hailong Sun, *Beihang University*; Xin Zhuang, *University of Toronto*; Nikola Grcevski, *Vena Solutions*; Ding Yuan, *University of Toronto*

## Potpourri

**EC-Cache: Load-Balanced, Low-Latency Cluster Caching with Online Erasure Coding** .....401  
K. V. Rashmi, *University of California, Berkeley*; Mosharaf Chowdhury and Jack Kosaian, *University of Michigan*; Ion Stoica and Kannan Ramchandran, *University of California, Berkeley*

**To Waffinity and Beyond: A Scalable Architecture for Incremental Parallelization of File System Code**...419  
Matthew Curtis-Maury, Vinay Devadas, Vania Fang, and Aditya Kulkarni, *NetApp, Inc.*

**CLARINET: WAN-Aware Optimization for Analytics Queries** .....435  
Raajay Viswanathan, *University of Wisconsin—Madison*; Ganesh Ananthanarayanan, *Microsoft*; Aditya Akella, *University of Wisconsin—Madison*

**JetStream: Cluster-Scale Parallelization of Information Flow Queries** .....451  
Andrew Quinn, David Devecsery, Peter M. Chen, and Jason Flinn, *University of Michigan*

## Fault Tolerance and Consensus

**Just say NO to Paxos Overhead: Replacing Consensus with Network Ordering** .....467  
Jialin Li, Ellis Michael, Naveen Kr. Sharma, Adriana Szekeres, and Dan R. K. Ports, *University of Washington*

**XFT: Practical Fault Tolerance beyond Crashes** .....485  
Shengyun Liu, *National University of Defense Technology*; Paolo Viotti, *EURECOM*; Christian Cachin, *IBM Research—Zurich*; Vivien Quéma, *Grenoble Institute of Technology*; Marko Vukolić, *IBM Research—Zurich*

**Realizing the Fault-Tolerance Promise of Cloud Storage Using Locks with Intent** .....501  
Srinath Setty, *Microsoft Research*; Chunzhi Su, *The University of Texas at Austin and Microsoft Research*; Jacob R. Lorch and Lidong Zhou, *Microsoft Research*; Hao Chen, *Shanghai Jiao Tong University and Microsoft Research*; Parveen Patel and Jinglei Ren, *Microsoft Research*

**Consolidating Concurrency Control and Consensus for Commits under Conflicts** .....517  
Shuai Mu and Lamont Nelson, *New York University*; Wyatt Lloyd, *University of Southern California*; Jinyang Li, *New York University*

## Friday, November 4, 2016

### Security

**Ryoan: A Distributed Sandbox for Untrusted Computation on Secret Data** .....533  
Tyler Hunt, Zhiting Zhu, Yuanzhong Xu, Simon Peter, and Emmett Witchel, *The University of Texas at Austin*

**Unobservable Communication over Fully Untrusted Infrastructure** .....551  
Sebastian Angel, *The University of Texas at Austin and New York University*; Srinath Setty, *Microsoft Research*

**Alpenhorn: Bootstrapping Secure Communication Without Leaking Metadata** .....571  
David Lazar and Nickolai Zeldovich, *MIT CSAIL*

**Big Data Analytics over Encrypted Datasets with Seabed** .....587  
Antonis Papadimitriou, *University of Pennsylvania and Microsoft Research India*; Ranjita Bhagwan, Nishanth Chandran, and Ramachandran Ramjee, *Microsoft Research India*; Andreas Haeberlen, *University of Pennsylvania*; Harmeet Singh and Abhishek Modi, *Microsoft Research India*; Saikrishna Badrinarayanan, *University of California, Los Angeles and Microsoft Research India*

(Friday, November 4, continues on next page)

## Troubleshooting

**Non-Intrusive Performance Profiling for Entire Software Stacks Based on the Flow Reconstruction Principle**.....603  
Xu Zhao, Kirk Rodrigues, Yu Luo, Ding Yuan, and Michael Stumm, *University of Toronto*

**Early Detection of Configuration Errors to Reduce Failure Damage**.....619  
Tianyin Xu, Xinxin Jin, Peng Huang, and Yuanyuan Zhou, *University of California, San Diego*; Shan Lu, *University of Chicago*; Long Jin, *University of California, San Diego*; Shankar Pasupathy, *NetApp, Inc.*

**Kraken: Leveraging Live Traffic Tests to Identify and Resolve Resource Utilization Bottlenecks in Large Scale Web Services** .....635  
Kaushik Veeraraghavan, Justin Meza, David Chou, Wonho Kim, Sonia Margulis, Scott Michelson, Rajesh Nishtala, Daniel Obenshain, Dmitri Perelman, and Yee Jiun Song, *Facebook Inc.*

## Operating Systems II

**CertiKOS: An Extensible Architecture for Building Certified Concurrent OS Kernels** .....653  
Ronghui Gu, Zhong Shao, Hao Chen, Xiongnan (Newman) Wu, Jieung Kim, Vilhelm Sjöberg, and David Costanzo, *Yale University*

**EbbRT: A Framework for Building Per-Application Library Operating Systems** .....671  
Dan Schatzberg, James Cadden, Han Dong, Orran Krieger, and Jonathan Appavoo, *Boston University*

**SCONE: Secure Linux Containers with Intel SGX**.....689  
Sergei Arnautov, Bohdan Trach, Franz Gregor, Thomas Knauth, and Andre Martin, *Technische Universität Dresden*; Christian Priebe, Joshua Lind, Divya Muthukumaran, Dan O’Keeffe, and Mark L Stillwell, *Imperial College London*; David Goltzsche, *Technische Universität Braunschweig*; Dave Eyers, *University of Otago*; Rüdiger Kapitza, *Technische Universität Braunschweig*; Peter Pietzuch, *Imperial College London*; Christof Fetzer, *Technische Universität Dresden*

**Coordinated and Efficient Huge Page Management with Ingens** .....705  
Youngjin Kwon, Hangchen Yu, and Simon Peter, *The University of Texas at Austin*; Christopher J. Rossbach, *The University of Texas at Austin and VMware*; Emmett Witchel, *The University of Texas at Austin*

## Cloud Systems II

**Diamond: Automating Data Management and Storage for Wide-Area, Reactive Applications**.....723  
Irene Zhang, Niel Lebeck, Pedro Fonseca, Brandon Holt, Raymond Cheng, Ariadna Norberg, Arvind Krishnamurthy, and Henry M. Levy, *University of Washington*

**Slicer: Auto-Sharding for Datacenter Applications** .....739  
Atul Adya, Daniel Myers, Jon Howell, Jeremy Elson, Colin Meek, Vishesh Khemani, Stefan Fulger, Pan Gu, Lakshminath Bhuvanagiri, Jason Hunter, and Roberto Peon, Larry Kai, Alexander Shraer, and Arif Merchant, *Google*; Kfir Lev-Ari, *Technion—Israel Institute of Technology*

**History-Based Harvesting of Spare Cycles and Storage in Large-Scale Datacenters** .....755  
Yunqi Zhang, *University of Michigan and Microsoft Research*; George Prekas, *École Polytechnique Fédérale de Lausanne (EPFL) and Microsoft Research*; Giovanni Matteo Fumarola and Marcus Fontoura, *Microsoft*; Íñigo Goiri and Ricardo Bianchini, *Microsoft Research*

**DQBarge: Improving Data-Quality Tradeoffs in Large-Scale Internet Services** .....771  
Michael Chow, *University of Michigan*; Kaushik Veeraraghavan, *Facebook, Inc.*; Michael Cafarella and Jason Flinn, *University of Michigan*