

19th USENIX Symposium on Operating Systems Design and Implementation (OSDI '25)

July 7–9, 2025
Boston, MA, USA

Monday, July 7

Distributed Systems and Data Centers I

Basilisk: Using Provenance Invariants to Automate Proofs of Undecidable Protocols	1
Tony Nuda Zhang and Keshav Singh, <i>University of Michigan</i> ; Tej Chajed, <i>University of Wisconsin-Madison</i> ; Manos Kapritsos, <i>University of Michigan</i> ; Bryan Parno, <i>Carnegie Mellon University</i>	
Deriving Semantic Checkers from Tests to Detect Silent Failures in Production Distributed Systems	19
Chang Lou, <i>University of Virginia</i> ; Dimas Shidqi Parikesit, <i>University of Virginia and Bandung Institute of Technology</i> ; Yujin Huang, <i>The Pennsylvania State University</i> ; Zhewen Yang and Senapati Diwangkara, <i>Johns Hopkins University</i> ; Yuzhuo Jing, <i>University of Michigan</i> ; Achmad Imam Kistijantoro, <i>Bandung Institute of Technology</i> ; Ding Yuan, <i>University of Toronto</i> ; Suman Nath, <i>Microsoft Research</i> ; Peng Huang, <i>University of Michigan</i>	
Picsou: Enabling Replicated State Machines to Communicate Efficiently	39
Reginald Frank, Micah Murray, Chawinphat Tankuranand, Junseo Yoo, Ethan Xu, and Natacha Crooks, <i>UC Berkeley</i> ; Suyash Gupta, <i>University of Oregon</i> ; Manos Kapritsos, <i>University of Michigan</i>	
FineMem: Breaking the Allocation Overhead vs. Memory Waste Dilemma in Fine-Grained Disaggregated Memory Management	57
Xiaoyang Wang and Yongkun Li, <i>University of Science and Technology of China</i> ; Kan Wu, <i>Google</i> ; Wenzhe Zhu and Yuqi Li, <i>University of Science and Technology of China</i> ; Yinlong Xu, <i>University of Science and Technology of China and Anhui Provincial Key Laboratory of High Performance Computing</i>	
To PRI or Not To PRI, That's the question	75
Yun Wang, <i>Shanghai Jiao Tong University</i> ; Liang Chen, Jie Ji, Xianting Tian, and Ben Luo, <i>Alibaba Group</i> ; Zhixiang Wei, Zhibai Huang, and Kailiang Xu, <i>Shanghai Jiao Tong University</i> ; Kaihuan Peng, Kaijie Guo, Ning Luo, Guangjian Wang, Shengdong Dai, Yibin Shen, and Jiesheng Wu, <i>Alibaba Group</i> ; Zhengwei Qi, <i>Shanghai Jiao Tong University</i>	
Enabling Efficient GPU Communication over Multiple NICs with FuseLink	91
Zhenghang Ren, Yuxuan Li, Zilong Wang, Xinyang Huang, Wenxue Li, Kaiqiang Xu, Xudong Liao, Yijun Sun, and Bowen Liu, <i>Hong Kong University of Science and Technology</i> ; Han Tian, <i>University of Science and Technology of China</i> ; Junxue Zhang, <i>Hong Kong University of Science and Technology</i> ; Mingfei Wang, <i>MetaX Integrated Circuits</i> ; Zhizhen Zhong, <i>Massachusetts Institute of Technology</i> ; Guyue Liu, <i>Peking University</i> ; Ying Zhang, <i>Meta</i> ; Kai Chen, <i>Hong Kong University of Science and Technology</i>	
Database Systems	
Tigon: A Distributed Database for a CXL Pod	109
Yibo Huang, Haowei Chen, and Newton Ni, <i>The University of Texas at Austin</i> ; Yan Sun, <i>University of Illinois Urbana-Champaign</i> ; Vijay Chidambaram, Dixin Tang, and Emmett Witchel, <i>The University of Texas at Austin</i>	
Mako: Speculative Distributed Transactions with Geo-Replication	129
Weihai Shen, <i>Stony Brook University</i> ; Yang Cui, <i>Google</i> ; Siddhartha Sen, <i>Microsoft Research</i> ; Sebastian Angel, <i>University of Pennsylvania</i> ; Shuai Mu, <i>Stony Brook University</i>	
Quake: Adaptive Indexing for Vector Search	153
Jason Mohoney, Devesh Sarda, and Mengze Tang, <i>University of Wisconsin-Madison</i> ; Shihabur Rahman Chowdhury and Anil Pacaci, <i>Apple</i> ; Ihab F. Ilyas, <i>University of Waterloo</i> ; Theodoros Rekatsinas, <i>Apple</i> ; Shivaram Venkataraman, <i>University of Wisconsin-Madison</i>	
Achieving Low-Latency Graph-Based Vector Search via Aligning Best-First Search Algorithm with SSD	171
Hao Guo and Youyou Lu, <i>Tsinghua University</i>	
Skybridge: Bounded Staleness for Distributed Caches	187
Robert Lyerly, <i>Meta Platforms Inc.</i> ; Scott Pruett, <i>unaffiliated</i> ; Kevin Doherty and Greg Rogers, <i>Meta Platforms Inc.</i> ; Nathan Bronson, <i>OpenAI</i> ; John Hugg, <i>Meta Platforms Inc.</i>	

AI + Systems I

KPerfIR: Towards a Open and Compiler-centric Ecosystem for GPU Kernel Performance Tooling on Modern AI Workloads	205
Yue Guan, <i>University of California, San Diego</i> ; Yuanwei Fang, <i>Meta</i> ; Keren Zhou, <i>George Mason University and OpenAI</i> ; Corbin Robeck and Manman Ren, <i>Meta</i> ; Zhongkai Yu, <i>University of California, San Diego</i> ; Yufei Ding, <i>University of California, San Diego, and Meta</i> ; Adnan Aziz, <i>Meta</i>	
Mirage: A Multi-Level Superoptimizer for Tensor Programs	221
Mengdi Wu and Xinhao Cheng, <i>Carnegie Mellon University</i> ; Shengyu Liu and Chunan Shi, <i>Peking University</i> ; Jianan Ji and Man Kit Ao, <i>Carnegie Mellon University</i> ; Praveen Velliengiri, <i>Pennsylvania State University</i> ; Xupeng Miao, <i>Purdue University</i> ; Oded Padon, <i>Weizmann Institute of Science</i> ; Zhihao Jia, <i>Carnegie Mellon University</i>	
QiMeng-Xpiler: Transcompiling Tensor Programs for Deep Learning Systems with a Neural-Symbolic Approach ..	239
Shouyang Dong, <i>University of Science and Technology of China, Cambricon Technologies, and Institute of Computing Technology, Chinese Academy of Sciences</i> ; Yuanbo Wen, Jun Bi, Di Huang, and Jiaming Guo, <i>Institute of Computing Technology, Chinese Academy of Sciences</i> ; Jianxing Xu and Ruibai Xu, <i>University of Science and Technology of China, Cambricon Technologies, and Institute of Computing Technology, Chinese Academy of Sciences</i> ; Xinkai Song and Yifan Hao, <i>Institute of Computing Technology, Chinese Academy of Sciences</i> ; Ling Li, <i>Institute of Software, Chinese Academy of Sciences, and University of Chinese Academy of Sciences</i> ; Xuehai Zhou, <i>University of Science and Technology of China</i> ; Tianshi Chen, <i>Cambricon Technologies</i> ; Qi Guo, <i>Institute of Computing Technology, Chinese Academy of Sciences</i> ; Yunji Chen, <i>Institute of Computing Technology, Chinese Academy of Sciences, and University of Chinese Academy of Sciences</i>	
WaferLLM: Large Language Model Inference at Wafer Scale	257
Congjie He, Yeqi Huang, and Pei Mu, <i>University of Edinburgh</i> ; Ziming Miao, Jilong Xue, Lingxiao Ma, and Fan Yang, <i>Microsoft Research</i> ; Luo Mai, <i>University of Edinburgh</i>	

Tuesday, July 8

AI + Systems II

BLITZSCALE: Fast and Live Large Model Autoscaling with $O(1)$ Host Caching	275
Dingyan Zhang, Haotian Wang, Yang Liu, and Xingda Wei, <i>Shanghai Jiao Tong University</i> ; Yizhou Shan, <i>Huawei Cloud</i> ; Rong Chen and Haibo Chen, <i>Shanghai Jiao Tong University</i>	
Bayesian Code Diffusion for Efficient Automatic Deep Learning Program Optimization	295
Isu Jeong and Seulki Lee, <i>Ulsan National Institute of Science and Technology</i>	
Training with Confidence: Catching Silent Errors in Deep Learning Training with Automated Proactive Checks ...	313
Yuxuan Jiang, Ziming Zhou, Boyu Xu, Beijie Liu, Runhui Xu, and Peng Huang, <i>University of Michigan</i>	
NEUTRINO: Fine-grained GPU Kernel Profiling via Programmable Probing	331
Songlin Huang and Chenshu Wu, <i>The University of Hong Kong</i>	
Principles and Methodologies for Serial Performance Optimization	357
Sujin Park, Mingyu Guan, Xiang Cheng, and Taesoo Kim, <i>Georgia Institute of Technology</i>	

Scheduling and Resource Management

Söße: One Network Telemetry Is All You Need for Per-flow Weighted Bandwidth Allocation at Scale	375
Weitao Wang and T. S. Eugene Ng, <i>Rice University</i>	
Decouple and Decompose: Scaling Resource Allocation with DeDe	393
Zhiying Xu and Minlan Yu, <i>Harvard University</i> ; Francis Y. Yan, <i>University of Illinois Urbana-Champaign</i>	
Quantum Virtual Machines	411
Runzhou Tao and Hongzheng Zhu, <i>University of Maryland, College Park</i> ; Jason Nieh, <i>Columbia University</i> ; Jianan Yao, <i>University of Toronto</i> ; Ronghui Gu, <i>Columbia University</i>	
QOS: Quantum Operating System	429
Emmanouil Giortamis, Francisco Romão, Nathaniel Tornow, and Pramod Bhatotia, <i>TU Munich</i>	
Scalio: Scaling up DPU-based JBOF Key-value Store with NVMe-oF Target Offload	449
Xun Sun, Mingxing Zhang, Yingdi Shan, Kang Chen, and Jinlei Jiang, <i>Tsinghua University</i> ; Yongwei Wu, <i>Tsinghua University, Quan Cheng Laboratory</i>	

Distributed Systems and Data Centers II

- Low End-to-End Latency atop a Speculative Shared Log with Fix-Ante Ordering** 465
Shreesha G. Bhat, Tony Hong, Xuhao Luo, Jiyu Hu, Aishwarya Ganesan, and Ramnatthan Alagappan, *University of Illinois Urbana-Champaign*
- Understanding Stragglers in Large Model Training Using What-if Analysis** 483
Jinkun Lin, *New York University*; Ziheng Jiang, Zuquan Song, Sida Zhao, and Menghan Yu, *ByteDance Seed*;
Zhanghan Wang, *New York University*; Chenyuan Wang, *ByteDance Seed*; Zuo Cheng Shi, *Zhejiang University*;
Xiang Shi, *ByteDance*; Wei Jia, Zherui Liu, Shuguang Wang, Haibin Lin, and Xin Liu, *ByteDance Seed*; Aurojit Panda and Jinyang Li, *New York University*
- Fork in the Road: Reflections and Optimizations for Cold Start Latency in Production Serverless Systems** 499
Xiaohu Chai, *Tsinghua University*; Ant Group; Tianyu Zhou, *Ant Group*; Keyang Hu, *Tsinghua University*; Jianfeng Tan, Tiwei Bie, Anqi Shen, Dawei Shen, Qi Xing, Shun Song, Tongkai Yang, Le Gao, Feng Yu, and Zhengyu He, *Ant Group*; Dong Du and Yubin Xia, *Shanghai Jiao Tong University*; Kang Chen, *Tsinghua University*; Yu Chen, *Quan Cheng Laboratory*; *Tsinghua University*
- Kamino: Efficient VM Allocation at Scale with Latency-Driven Cache-Aware Scheduling** 519
David Domingo, *Rutgers University*; Hugo Barbalho and Marco Molinaro, *Microsoft Research*; Kuan Liu, Abhisek Pan, David Dion, and Thomas Moscibroda, *Microsoft Azure*; Sudarsun Kannan, *Rutgers University*; Ishai Menache, *Microsoft Research*
- ZEN: Empowering Distributed Training with Sparsity-driven Data Synchronization** 537
Zhuang Wang, *Rice University*; Zhaozhuo Xu, *Stevens Institute of Technology*; Jingyi Xi, *unaffiliated*; Yuke Wang, Anshumali Shrivastava, and T. S. Eugene Ng, *Rice University*

Kernel and Operating Systems I

- Extending Applications Safely and Efficiently** 557
Yusheng Zheng, *UC Santa Cruz*; Tong Yu, *eunomia-bpf Community*; Yiwei Yang, *UC Santa Cruz*; Yanpeng Hu, *ShanghaiTech University*; Xiaozheng Lai, *South China University of Technology*; Dan Williams, *Virginia Tech*; Andi Quinn, *UC Santa Cruz*
- Tintin: A Unified Hardware Performance Profiling Infrastructure to Uncover and Manage Uncertainty** 575
Ao Li, Marion Sudvarg, Zihan Li, Sanjoy Baruah, Chris Gill, and Ning Zhang, *Washington University in St. Louis*
- Building Bridges: Safe Interactions with Foreign Languages through Omniglot** 595
Leon Schuermann and Jack Toubes, *Princeton University*; Tyler Potyondy and Pat Pannuto, *University of California San Diego*; Mae Milano and Amit Levy, *Princeton University*
- KRR: Efficient and Scalable Kernel Record Replay** 615
Tianren Zhang, *SmartX*; Sishuai Gong and Pedro Fonseca, *Purdue University*
- Deterministic Client: Enforcing Determinism on Untrusted Machine Code** 633
Zachary Yedidia, *Stanford University*; Geoffrey Ramseyer, *Stanford University and Stellar Development Foundation*; David Mazieres, *Stanford University*

Wednesday, July 9

Kernel and Operating Systems II

- Disentangling the Dual Role of NIC Receive Rings** 651
Boris Pismenny, *EPFL & NVIDIA*; Adam Morrison, *Tel Aviv University*; Dan Tsafir, *Technion – Israel Institute of Technology*
- XSched: Preemptive Scheduling for Diverse XPU** 671
Weihang Shen, Mingcong Han, Jialong Liu, Rong Chen, and Haibo Chen, *Institute of Parallel and Distributed Systems, Shanghai Jiao Tong University*
- OS Rendering Service Made Parallel with Out-of-Order Execution and In-Order Commit** 693
Yuanpei Wu and Dong Du, *Institute of Parallel and Distributed Systems, Shanghai Jiao Tong University*; *Engineering Research Center for Domain-specific Operating Systems, Ministry of Education*; Chao Xu, *Fields Lab, Huawei Central Software Institute*; Yubin Xia and Yang Yu, *Institute of Parallel and Distributed Systems, Shanghai Jiao Tong University*; *Engineering Research Center for Domain-specific Operating Systems, Ministry of Education*; Ming Fu, *Fields Lab, Huawei Central Software Institute*; Binyu Zang and Haibo Chen, *Institute of Parallel and Distributed Systems, Shanghai Jiao Tong University*; *Engineering Research Center for Domain-specific Operating Systems, Ministry of Education*

EMT: An OS Framework for New Memory Translation Architectures 711
Siyuan Chai, Jiyuan Zhang, Jongyul Kim, Alan Wang, Fan Chung, and Jovan Stojkovic, *University of Illinois Urbana-Champaign*; Weiwei Jia, *University of Rhode Island*; Dimitrios Skarlatos, *Carnegie Mellon University*; Josep Torrellas and Tianyin Xu, *University of Illinois Urbana-Champaign*

Tiered Memory Management Beyond Hotness 731
Jinshu Liu, Hamid Hadian, Hanchen Xu, and Huaicheng Li, *Virginia Tech*

AI + Systems III

NanoFlow: Towards Optimal Large Language Model Serving Throughput 749
Kan Zhu, *University of Washington*; Yufei Gao, *Tsinghua University and University of Washington*; Yilong Zhao, *University of Washington and University of California, Berkeley*; Liangyu Zhao, *University of Washington*; Gefei Zuo, *University of Michigan*; Yile Gu and Dedong Xie, *University of Washington*; Tian Tang and Qinyu Xu, *Tsinghua University and University of Washington*; Zihao Ye, Keisuke Kamahori, and Chien-Yu Lin, *University of Washington*; Ziren Wang, *Tsinghua University and University of Washington*; Stephanie Wang, Arvind Krishnamurthy, and Baris Kasikci, *University of Washington*

PipeThreader: Software-Defined Pipelining for Efficient DNN Execution 767
Yu Cheng, Lei Wang, and Yining Shi, *School of Computer Science, Peking University*; Yuqing Xia, Lingxiao Ma, Jilong Xue, and Yang Wang, *Microsoft Research*; Zhiwen Mo, *Imperial College London and Microsoft Research*; Feiyang Chen, *Shanghai Jiao Tong University and Microsoft Research*; Fan Yang and Mao Yang, *Microsoft Research*; Zhi Yang, *School of Computer Science, Peking University*

WLB-LLM: Workload-Balanced 4D Parallelism for Large Language Model Training 785
Zheng Wang, *University of California, San Diego, and Meta*; Anna Cai and Xinfeng Xie, *Meta*; Zaifeng Pan and Yue Guan, *University of California, San Diego*; Weiwei Chu, Jie Wang, Shikai Li, Jianyu Huang, Chris Cai, and Yuchen Hao, *Meta*; Yufei Ding, *University of California, San Diego, and Meta*

DecDEC: A Systems Approach to Advancing Low-Bit LLM Quantization 803
Yeonhong Park, Jake Hyun, Hojoon Kim, and Jae W. Lee, *Seoul National University*

File and Storage Systems

Stripeless Data Placement for Erasure-Coded In-Memory Storage 821
Jian Gao, Jiwu Shu, Bin Yan, and Yuhao Zhang, *Tsinghua University*; Keji Huang, *Huawei Technologies Co., Ltd*

POWER Never Corrupts: Tool-Agnostic Verification of Crash Consistency and Corruption Detection 839
Hayley LeBlanc, *University of Texas at Austin*; Jacob R. Lorch and Chris Hawblitzel, *Microsoft Research*; Cheng Huang and Yiheng Tao, *Microsoft*; Nikolai Zeldovich, *MIT CSAIL and Microsoft Research*; Vijay Chidambaram, *University of Texas at Austin*

Fast and Synchronous Crash Consistency with Metadata Write-Once File System 859
Yanqi Pan, Wen Xia, Yifeng Zhang, Xiangyu Zou, and Hao Huang, *Harbin institute of Technology, Shenzhen*; Zhenhua Li, *Tsinghua University*; Chentao Wu, *Shanghai Jiao Tong University*

Decentralized, Epoch-based F2FS Journaling with Fine-grained Crash Recovery 879
Yaotian Cui and Zhiqi Wang, *The Chinese University of Hong Kong, China*; Renhai Chen, *College of Intelligence and Computing, Tianjin University, China*; Zili Shao, *The Chinese University of Hong Kong, China*

Okapi: Decoupling Data Striping and Redundancy Grouping in Cluster File Systems 897
Sanjith Athlur and Timothy Kim, *Carnegie Mellon University*; Saurabh Kadekodi, *Google*; Francisco Maturana and Xavier Ramos, *Carnegie Mellon University*; Arif Merchant, *Google*; K. V. Rashmi and Gregory R. Ganger, *Carnegie Mellon University*

Privacy and Security

- Compass: Encrypted Semantic Search with High Accuracy** 915
Jinhao Zhu, *UC Berkeley*; Liana Patel, *Stanford University*; Matei Zaharia and Raluca Ada Popa, *UC Berkeley*
- Weave: Efficient and Expressive Oblivious Analytics at Scale** 939
Mahdi Soleimani, Grace Jia, and Anurag Khandelwal, *Yale University*
- Paralegal: Practical Static Analysis for Privacy Bugs.** 957
Justus Adam, Carolyn Zech, Livia Zhu, Sreshtaa Rajesh, Nathan Harbison, Mithi Jethwa, Will Crichton, Shriram Krishnamurthi, and Malte Schwarzkopf, *Brown University*
- MettEagle: Costs and Benefits of Implementing Containers on Microkernels** 979
Till Miemietz, Viktor Reusch, and Matthias Hille, *Barkhausen Institut*; Lars Wrenger, *Leibniz-Universität Hannover*; Jana Eisoldt, *Barkhausen Institut*; Jan Klötzke, *Kernkonzept GmbH*; Max Kurze, *Technische Universität Dresden*; Adam Lackorzynski, *Technische Universität Dresden and Kernkonzept GmbH*; Michael Roitzsch, *Barkhausen Institut*; Hermann Härtig, *Barkhausen Institut and Technische Universität Dresden*