

# **2025 IEEE International Conference on Cluster Computing (CLUSTER 2025)**

**Edinburgh, United Kingdom  
2-5 September 2025**



**IEEE Catalog Number: CFP25235-POD  
ISBN: 979-8-3315-3020-4**

**Copyright © 2025 by the Institute of Electrical and Electronics Engineers, Inc.  
All Rights Reserved**

*Copyright and Reprint Permissions:* Abstracting is permitted with credit to the source. Libraries are permitted to photocopy beyond the limit of U.S. copyright law for private use of patrons those articles in this volume that carry a code at the bottom of the first page, provided the per-copy fee indicated in the code is paid through Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923.

For other copying, reprint or republication permission, write to IEEE Copyrights Manager, IEEE Service Center, 445 Hoes Lane, Piscataway, NJ 08854. All rights reserved.

***\*\*\* This is a print representation of what appears in the IEEE Digital Library. Some format issues inherent in the e-media version may also appear in this print version.***

IEEE Catalog Number:	CFP25235-POD
ISBN (Print-On-Demand):	979-8-3315-3020-4
ISBN (Online):	979-8-3315-3019-8
ISSN:	1552-5244

**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  
E-mail: [curran@proceedings.com](mailto:curran@proceedings.com)  
Web: [www.proceedings.com](http://www.proceedings.com)

CURRAN ASSOCIATES INC.  
**proceedings**  
.com

# Contents

<b>Welcome Message from the General Chairs</b>	<b>4</b>
<b>Welcome Message from the Technical Program Chairs</b>	<b>5</b>
<b>Committees</b>	<b>7</b>
<b>Papers</b>	<b>13</b>
RAN: Accelerating Data Repair with Available Nodes in Erasure-Coded Storage . . . . .	13
Scalable and Fast Inference Serving via Hybrid Communication Scheduling on Heterogeneous Networks . . . . .	25
Multi-agent Independent PPO-based Automatic ECN Tuning for High-Speed Data Center Networks . . . . .	37
SplitQuant: Resource-Efficient LLM Offline Serving on Heterogeneous GPUs via Phase-Aware Model Partition and Adaptive Quantization	48
Rock: Serving Multimodal Models in Cloud with Heterogeneous-Aware Resource Orchestration for Thousands of LoRA Adapters . . . . .	59
Accelerating Key-Value Data Structures Using AVX-512 SIMD Extensions . . . . .	72
EquilibRIO: Taming the I/O Tides in High-Performance Computing . . . . .	84
FIFO-MEP: An Efficient Multi-Eviction-Point FIFO Cache with Stable Demotion for Burst-Oriented Access Mitigation . . . . .	96
Are We There Yet? Predicting the Queue Wait Times for HPC Jobs . . . . .	109
GreenK8s: Green-aware Scheduling for Sustainable Kubernetes Cluster Management . . . . .	121
Revisiting Fragmentation for Deduplication in Clustered Primary Storage Systems . . . . .	133
PIAR: Path-Improved Adaptive Routing for Dragonfly Networks . . . . .	145
A Versatile Simulated Data Transport Layer for In Situ Workflows Performance Evaluation . . . . .	156
Lessons from Profiling and Optimizing Placement in AMR Codes . . . . .	167
CFseq: A Framework for Constructing Compression-Friendly Field Sequences for Network Logs . . . . .	181
Towards High-Performance and Portable Molecular Docking on CPUs through Vectorization . . . . .	193
SoCL: Scalable and Latency-Optimized Microservices in Serverless Edge Computing . . . . .	207
Proactive SSD Failure Prediction with A Gradient-Guided LSTM-xLSTM Hybrid Model . . . . .	218
DDRM: An SLO-aware Deep Dynamic Resource Management Framework for Microservices . . . . .	229
Parallel tall-and-skinny QR factorization based on LU-CholeskyQR algorithm . . . . .	241
PRT: An Efficient Pipeline Reuse Technology for Large Models Training	251
NSYS2PRV: detailed and quantitative analysis of large-scale GPU execution traces with Paraver . . . . .	262
A Pattern-Aware Finite Element Matrix Assembly Method on GPUs . . . . .	274

Communication Notification through User-Level Interrupts for the BXI Network . . . . .	284
DaCe AD: Unifying High-Performance Automatic Differentiation for Machine Learning and Scientific Computing . . . . .	294
Scaling Deep Learning Molecular Dynamics to 500M Atoms on 4096-Node ARMv8 Clusters . . . . .	307
Closing the HPC-Cloud Convergence Gap: Multi-Tenant Slingshot RDMA for Kubernetes . . . . .	319
UNICONN: A Uniform High-Level Communication Library for Portable Multi-GPU Programming . . . . .	329
TRACE: A Targeted Recommender for VM Assignment in Cloud Environment . . . . .	341
Deadline-Aware Resource Allocation and Scheduling of Serverless Workloads on Heterogeneous Clusters . . . . .	352
Parallel Selected Inversion of Block-Tridiagonal with Arrowhead Matrices	363
Bridging Metadata Service and CXL: A Metadata-grained and Directory-aware Storage Engine for Distributed Storage Systems . . . . .	375
Efficient Multi-GPU Programming in Python: Reducing Synchronization and Access Overheads . . . . .	387
Cache Less to Save More: A Cost-Based Distributed Caching Strategy for ICN . . . . .	397
BMPipe: Bubble-Memory Co-optimization Strategy Planner for Very-large DNN Training . . . . .	409
Fine-grain energy consumption modeling of HPC task-based programs	421
Detecting Silent Data Corruption From Hardware Counters . . . . .	433
Towards Dynamic Message Passing Protocols for Stencil-based Communication Patterns . . . . .	446
Cascade: a Collaborative Algorithm for Scalable And Efficient Neighborhood Allgather . . . . .	458
Capricorn: Efficient In-Memory Checkpointing for MoE Model Training with Dynamicity Awareness . . . . .	471