

Mustafa Ozan Karsavuran

Postdoctoral Scholar, Lawrence Berkeley National Laboratory, Berkeley, California, 94720 US

✉ ozan.karsavuran@gmail.com | ✉ MOKarsavuran@lbl.gov | 📄 ozan-karsavuran-421a7419

SUMMARY

Experienced researcher specializing in high-performance computing, parallel algorithms, and sparse matrix/tensor computations. Skilled in MPI and OpenMP programming, GPU acceleration, and hypergraph partitioning models. Proven track record of performance optimization, algorithm development, and successful collaboration on large-scale scientific computing projects.

EDUCATION

Bilkent University, Dept. of Computer Engineering, Turkey

- **Doctor of Philosophy** GPA: 3.47/4.00 September 2014 – August 2020
Thesis: Reducing Communication Overhead in Sparse Matrix and Tensor Computations
Advisor: Prof. Dr. Cevdet Aykanat
- **Master of Science** GPA: 3.60/4.00 September 2012 – September 2014
Thesis: Increasing Data Reuse in Parallel SpMV and SpMTV Multiply on Shared-Memory Architectures
Advisor: Prof. Dr. Cevdet Aykanat
- **Bachelor of Science** GPA: 3.44/4.00 September 2007 – June 2012

RESEARCH EXPERIENCE

Postdoctoral Scholar

Lawrence Berkeley National Laboratory

November 2022 – *present*
Berkeley, California, U.S.

★ Sparse Symmetric Matrix Factorization

- Right-looking sparse $A = LDL^T$ factorization 2023
 - ◇ Developed efficient pivoting strategies for factoring sparse general symmetric matrices
 - ◇ Achieved higher accuracy and reduced fill-in
- Right-looking sparse Cholesky factorization 2023
 - ◇ Accelerated a serial Cholesky algorithm on GPU and achieved $4\times$ speedup and published the results in SC24-W
 - ◇ Developed a serial right-looking supernodal sparse Cholesky factorization
 - ◇ Obtained a Cholesky algorithm which needs no temporary working storage and avoids assembly operations

Postdoctoral Researcher

Bilkent University

September 2020 – October 2022
Ankara, Turkey

★ Stochastic Gradient Descent (SGD) for Matrix Completion

- Scaling stratified SGD for distributed matrix completion 2022
 - ◇ Collaborated on point-to-point (P2P) communication scheme and related hold-and-combine algorithm
 - ◇ Implemented a computational load balancing method and an HP model which minimize the communication volume
 - ◇ Obtained about up to $15\times$ faster parallel runtime and published the results in IEEE TKDE
- Reducing stale data usage and bandwidth requirement in SGD 2022
 - ◇ Contribute on the design of the MPI-based SGD with sub-iterations
 - ◇ Design and implement an HP model minimizing both staleness and bandwidth requirement in SGD
 - ◇ Obtained up to 34% reduction in parallel runtime and published the results in IEEE TC
- Hybrid Parallel SGD 2022
 - ◇ Contribute on the design and implementation of the MPI+POSIX threads based hybrid parallel SGD
 - ◇ Obtained up to $6\times$ better throughput and published the results in KNOSYS

★ Parallel Sparse Tensor Decomposition

- Hiding latency of the sparse P2P communications into dense all reduce communications 2021
 - ◇ Collaborated on reorganizing the CPD algorithm for embedding P2P messages into ALLREDUCE messages
 - ◇ Designed and implemented an HP model which minimizes the concurrent communication volume in the embedded ALLREDUCE
 - ◇ Scaled CPD on up to 4096 processors and published the results in IEEE TPDS
- General medium-grain sparse tensor partitioning for distributed CPD 2019
 - ◇ Designed and implemented an HP model for medium-grain partitioning without any topological constraint
 - ◇ Utilized RB paradigm to boost performance at each level of the partitioning
 - ◇ Conducted experiments with 10 real-world tensors on 1024 cores and published the results in IEEE TPDS
- Locality-aware fiber and slice reordering for shared-memory MTTKRP 2017
 - ◇ Implemented an HP model for reordering fibers and/or slices of tensors for increasing cache locality during MTTKRP
 - ◇ Adopted SPLATT's OpenMP based MTTKRP and conduct experiments

★ Large Scale Benchmarking

- Code owner of the NEMO package in the PRACE-6IP T7.4.A activity 2021
 - ◇ Prepared architecture specific build files for both NEMO and XIOS packages for six tier-0 HPC systems
 - ◇ Benchmarked the NEMO and XIOS on HAWK, TGCC Joliot Curie, JUWELS, Marconi100, MareNostrum and SuperMUC-NG using up to 10,000 cores
 - ◇ Contribute to the deliverable PRACE-6IP-D7.4: Evaluation of Benchmark Performance

★ Other Hypergraph Partitioning Models

- Simultaneous computational and data load balancing of the processors on distributed-memory setting 2022
 - ◇ Collaborated on design of two-constraint HP models which encodes computational and data load simultaneously
 - ◇ Collaborated on experiments with two different applications and published the results in SIAM J. Sci. Comput.

Teaching and Research Assistant

Bilkent University

September 2012 – June 2020

Ankara, Turkey

★ Sparse Matrix Vector Multiplication (SpMV) and Sparse Matrix Dense Matrix Multiplication (SpMM)

- Volume balancing and latency minimization for reduce operations 2018
 - ◇ Formulate a novel vertex weighting scheme for the HP model which balance volume loads of processors
 - ◇ Implement a refinement algorithm called during the RB and decrease the increase in the communication volume
 - ◇ Perform extensive experiments on 512 processors for 70 matrices
 - ◇ Obtain 30% faster parallel runtime in column-parallel SpMV and published in IEEE TPDS
- Locality-aware shared-memory parallel $y \leftarrow AA^T x$ on many-core processors 2014
 - ◇ Implement OpenMP-based $y \leftarrow AA^T x$ which achieve reusing A-matrix nonzeros and vector entries
 - ◇ Conduct detailed experiments on Intel Xeon Phi co-processor running in offload mode
 - ◇ Obtain 20% reduction in parallel runtime and published in IEEE TPDS

★ Generalized Sparse Matrix Matrix Multiplication (SpGEMM)

- Efficient Vectorization of SpGEMM 2016
 - ◇ Transform an $C = ADB$ instance into $C = Zd$ SpMV instance by multiplying A-matrix with B-matrix columns
 - ◇ Implement efficiently vectorized OpenMP based $C = Zd$ using AVX instructions
 - ◇ Conduct experiments on Intel Xeon Phi co-processor and Xeon processor

TECHNICAL SKILLS

Advanced in: C, C++, MPI, OpenMP

Familiar with: CUDA, Fortran, Python, Java SE, C#, MATLAB, Assembly (MIPS and Intel 8051), PHP, SQL

JOURNAL PUBLICATIONS

- K. Büyükkaya, **M. O. Karsavuran** and C. Aykanat, “Stochastic Gradient Descent for Matrix Completion: Hybrid Parallelization on Shared- and Distributed-Memory Systems” in Knowledge-based Systems, vol. 283, pp. 111176, Jan. 2024. doi: 10.1016/j.knosys.2023.111176
- N. Abubaker, O. Çağlayan, **M. O. Karsavuran** and C. Aykanat, “Minimizing Staleness and Communication Overhead in Distributed SGD for Collaborative Filtering” in IEEE Transactions on Computers, vol. 72, no. 10, pp. 2925-2937, Oct. 2023. doi: 10.1109/TC.2023.3275107
- N. Abubaker, **M. O. Karsavuran** and C. Aykanat, “Scaling Stratified Stochastic Gradient Descent for Distributed Matrix Completion,” in IEEE Transactions on Knowledge and Data Engineering, vol. 35, no. 10, pp. 10603-10615, Oct. 2023. doi: 10.1109/TKDE.2023.3253791
- M.F. Çeliktug, **M. O. Karsavuran**, S. Acer and C. Aykanat, “Simultaneous Computational and Data Load Balancing in Distributed-Memory Setting,” in SIAM Journal on Scientific Computing, vol. 44, no. 6, pp. C399-C424, Nov. 2022. doi: 10.1137/22M1485772
- N. Abubaker, **M. O. Karsavuran** and C. Aykanat, “Scalable Unsupervised ML: Latency Hiding in Distributed Sparse Tensor Decomposition,” in IEEE Transactions on Parallel and Distributed Systems, vol. 33, no. 11, pp. 3028-3040, Nov. 2022. doi: 10.1109/TPDS.2021.3128827
- **M. O. Karsavuran**, S. Acer and C. Aykanat, “Partitioning Models for General Medium-Grain Parallel Sparse Tensor Decomposition,” in IEEE Transactions on Parallel and Distributed Systems, vol. 32, no. 1, pp. 147-159, Jan. 2021. doi: 10.1109/TPDS.2020.3012624
- **M. O. Karsavuran**, S. Acer and C. Aykanat, “Reduce Operations: Send Volume Balancing While Minimizing Latency,” in IEEE Transactions on Parallel and Distributed Systems, vol. 31, no. 6, pp. 1461-1473, June 2020. doi: 10.1109/TPDS.2020.2964536
- **M. O. Karsavuran**, K. Akbudak and C. Aykanat, “Locality-Aware Parallel Sparse Matrix-Vector and Matrix-Transpose-Vector Multiplication on Many-Core Processors,” in IEEE Transactions on Parallel and Distributed Systems, vol. 27, no. 6, pp. 1713-1726, June 2016. doi: 10.1109/TPDS.2015.2453970

CONFERENCE PUBLICATIONS

- **M. O. Karsavuran**, E. G. Ng, B. W. Peyton, “GPU Accelerated Sparse Cholesky Factorization” in SC24-W: Workshops of the International Conference for High Performance Computing, Networking, Storage and Analysis, Atlanta, GA, USA, 2024, pp. 703-707. doi: 10.1109/SCW63240.2024.00098

TALKS & POSTERS

- Sparse Cholesky Factorization Utilizing GPUs, SIAM CSE25, TX, USA.
- Heuristics for Robust Factorization of Sparse Symmetric Indefinite Matrices, SIAM LA24, Paris, France.
- Sparse Tensor Partitioning for Scalable Distributed CPD-ALS, SIAM PP22, Virtual.
- Medium-Grain Partitioning for Sparse Tensor Decomposition, SIAM CSE21, Virtual.
- Exploiting Matrix Reuse and Data Locality in Sparse Matrix-Vector and Matrix-Transpose-Vector Multiplication on Many-Core Architectures, SIAM CSC16, NM, USA.

PROFESSIONAL SERVICE

- **Reviewer** in ACM Transactions on Architecture and Code Optimization
 - **Reviewer** in The Journal of Supercomputing
 - **Reviewer** in CCPE (Concurrency and Computation: Practice and Experience)
 - **Reviewer** in PPAM22 (14th International Conference on Parallel Processing and Applied Mathematics)
 - **Reviewer** in BAŞARIM 2020 (6. Ulusal Yüksek Başarımlı Hesaplama Konferansı / 6th National Conference on High Performance Computing)
 - **Reviewer** in IPDPS 2018 (31st IEEE International Parallel & Distributed Processing Symposium)
-