

Gennady Pekhimenko

CONTACT INFORMATION	University of Toronto Computer Science Dept. (BA 5232) 40 St. George Street Toronto ON, M5S2E4	Work: (+1) 416-946-0250 Mobile: (+1) 647-916-6900 E-mail: pekhimenko@cs.toronto.edu Webpage: www.cs.toronto.edu/~pekhimenko/
RESEARCH INTERESTS	Systems, Computer Architecture, Applied Machine Learning	
EDUCATION	Carnegie Mellon University, USA <i>PhD in Computer Science, Computer Science Dept.</i> July 2016 Thesis: “Practical Data Compression for Modern Memory Hierarchies” Advisors: Todd C. Mowry and Onur Mutlu University of Toronto, Canada <i>MS in Computer Science</i> Department of Computer Science Jan 2008 Thesis: “Machine Learning Algorithms for Choosing Compiler Heuristics” Advisor: Angela Demke Brown Moscow State University, Russia <i>Diploma (5-year program) in Applied Mathematics & Computer Science</i> Faculty of Computational Mathematics and Cybernetics, Department of System Programming May 2004 Thesis: “Performance Analysis of MPI-Programs” Advisor: Victor A. Krukov	
APPOINTMENT	Vector Institute, Canada <i>Faculty Member, CIFAR AI Chair, Vector Institute</i> Aug 2019 – present <i>Faculty Affiliate, Vector Institute</i> May 2018 – Aug 2019 University of Toronto, Canada <i>Assistant Professor, Computer Science Department</i> June 2017 – present <i>Assistant Professor, Electrical & Computer Engineering Dept.</i> Jan 2018 – present	
AWARDS & HONORS	♦ ISCA Hall of Fame June 2021 at least 8 ISCA papers as a co-author ♦ IEEE MICRO Top Picks 2020 – 2021 A list of top papers from all computer architecture conferences that year ♦ HiPEAC 2020 Paper Award 2020 – 2020 Awarded for MLPerf Inference paper at ISCA’20 ♦ Amazon AWS Machine Learning Research Award 2020 – 2021 \$ 40,000 USD award and \$ 80,000 USD in AWS cloud credits. ♦ Facebook Faculty Research Award 2020 – 2021 \$ 49,500 USD award. One of the nine winners based on 132 proposals from 74 schools ♦ IEEE MICRO Top Picks Honorable Mention 2019 – 2020 A list of top papers from all computer architecture conferences that year ♦ CIFAR AI Chair 2019 – 2024 \$1,000,000 CAD award ♦ Connaught New Researcher Award 2018 – 2020 \$ 10,000 CAD award	

◇ NVIDIA Graduate Fellowship 5 winners nation-wide	2015 – 2016
◇ First place in ACM SRC (Student Research Competition) Energy-Efficient Data Compression for GPU Memory Systems @ ASPLOS'15	Mar 2015
◇ Qualcomm Innovation Fellowship Finalist Together with Nandita Vijaykumar. Selected as one of 35 out of 146 teams	2015 – 2016
◇ Facebook Fellowship Finalist \$500 cash prize	2015 – 2016
◇ Microsoft Research Fellowship 12 winners nation-wide	2013 – 2015
◇ Qualcomm Innovation Fellowship Together with Chris Fallin. 10 winner teams nation-wide	2013 – 2014
◇ First Heidelberg Laureate Forum Invitation Young researcher of the US delegation	Sep 2013
◇ Second place in ACM SRC (Student Research Competition) Linearly Compressed Pages: A Main Memory Compression Framework with Low Complexity and Low Latency @ PACT'12	Sep 2012
◇ Alexander Graham Bell Canada Graduate Scholarship NSERC (Canada's NSF) CGS-D2 Scholarship	2012 – 2013
◇ IBM First Patent Application Award Achievement \$2000 cash prize	Jan 2010
◇ Wolfond Scholarship University of Toronto Scholarship for high academic achievements	2006 – 2007
◇ Best Student Award Selected as the best student in MSU, CS Department	Apr 2003

PUBLICATIONS AND CITATIONS OVERVIEW Published 39 peer-reviewed conference papers, 10 journal papers/book chapters, and 4 peer-reviewed workshop papers in venues such as ISCA, MICRO, ASPLOS, HPCA, MLSys, UIST, USENIX ATC, PACT, SIGMETRICS, IISWC, TACO, CAL, Design & Test, and Oxford Bioinformatics. Awarded 3 patents. Given 67 conference/workshop/invited talks.

The total number of citations: **3255**, h-index: **28**
(based on Google Scholar information on September 11, 2021).

PEER-REVIEWED CONFERENCE PUBLICATIONS	<p>39. Omar Mohamed Awad, Mostafa Mahmoud, Isak Edo, Ali Hadi Zadeh, Ciaran Bannon, Anand Jayarajan, Gennady Pekhimenko, Andreas Moshovos. <i>FPRaker: A Processing Element For Accelerating Neural Network Training</i>. International Symposium on Microarchitecture (MICRO'21). October 2021.</p> <p>38. Geoffrey Yu, Pavel Golikov, YuBo Gao, Gennady Pekhimenko. <i>Habitat: Prediction-guided Hardware Selection for Deep Neural Network Training</i>. USENIX Annual Technical Conference (ATC'21). July 2021.</p> <p>37. Ziqi Wang, Michael A. Kozuch, Todd C. Mowry, Vivek Seshadri, Chulhwan Choo, Gennady Pekhimenko. <i>NVOverlay: Efficient, Scalable and Flexible Full-System Checkpointing on NVM with Low Write Amplification</i>. ACM/IEEE International Symposium on Computer Architecture (ISCA'21). June 2021.</p> <p>36. Shang Wang, Peiming Yang, Yuxuang Zheng, Xin Li, Gennady Pekhimenko. <i>Horizontally Fused Training Array: An Effective Hardware Utilization Squeezer for Training Novel Deep Learning Models</i>. Machine Learning and Systems Conference (MLSys'21). April 2021.</p> <p>35. James Gleeson, Srivatsan Krishnan, Moshe Gabel, Vijay Janapa Reddi, Eyal de Lara, Gennady Pekhimenko.</p>
--	--

- RL-Scope: Cross-stack Profiling for Deep Reinforcement Learning Workloads*. Machine Learning and Systems Conference (**MLSys'21**). April 2021.
34. Yaoyao Ding, Ligeng Zhu, Zhihao Jia, Gennady Pekhimenko, Song Han.
IOS: An Inter-Operator Scheduler for CNN Acceleration. Machine Learning and Systems Conference (**MLSys'21**). April 2021.
 33. Isak Edo Vivancos, Sayeh Sharify, Milos Nikolic, Ciaran Bannon, Mostafa Mahmoud, Alberto Delmás Lascorz, Gennady Pekhimenko, Andreas Moshovos.
Boveda: Building an On-Chip Deep Learning Memory Hierarchy Brick by Brick. Machine Learning and Systems Conference (**MLSys'21**). April 2021.
 32. Anand Jayarajan, Kimberly Hau, Andrew Goodwin, Gennady Pekhimenko.
LifeStream: A High-performance Stream Processing Engine for Periodic Streams. International Conference on Architectural Support for Programming Languages and Operating Systems (**ASPLOS'21**).
 31. Mostafa Mahmoud, Isak Edo Vivancos, Ali Hadi Zadeh, Omar Mohamed Awad, Jorge Albericio, Gennady Pekhimenko, Andreas Moshovos.
TensorDash: Exploiting Sparsity to Accelerate Deep Neural Network Training. International Symposium on Microarchitecture (**MICRO'20**). October 2020.
 30. Geoffrey Yu, Tovi Grossman, Gennady Pekhimenko.
Skyline: Interactive In-editor Computational Performance Profiling for Deep Neural Network Training. ACM Symposium on User Interface Software and Technology (**UIST'20**). October 2020.
 29. Hongyu Zhu, Amar Phanishayee, Gennady Pekhimenko.
Daydream: Accurately Estimating the Efficacy of Performance Optimizations for DNN Training. USENIX Annual Technical Conference (**ATC'20**). July 2020.
 28. Bojian Zheng, Nandita Vijaykumar, Gennady Pekhimenko.
Echo: Compiler-based GPU Memory Footprint Reduction for LSTM RNN Training. ACM/IEEE International Symposium on Computer Architecture (**ISCA'20**). June 2020.
 27. Vijay Janapa Reddi, Christine Cheng, David Kanter, Peter Mattson, Guenther Schmuelling, Carole-Jean Wu, Brian Anderson, Maximilien Breughe, Mark Charlebois, William Chou, Ramesh Chukka, Cody Coleman, Sam Davis, Pan Deng, Greg Diamos, Jared Duke, Dave Fick, J. Scott Gardner, Itay Hubara, Sachin Idgunji, Thomas B. Jablin, Jeff Jiao, Tom St. John, Pankaj Kanwar, David Leei, Jeffery Liao, Anton Lokhmotov, Francisco Massa, Peng Meng, Paulius Micikevicius, Colin Osborne, Gennady Pekhimenko, Arun Tejus, Raghunath Rajan, Dilip Sequeira, Ashish Sirasao, Fei Suni, Hanlin Tang, Michael Thomson, Frank Wei, Ephrem Wu, Lingjie Xu, Koichi Yamada, Bing Yu, George Yuan, Aaron Zhong, Peizhao Zhang, Yuchen Zhou.
MLPerf Inference Benchmark. ACM/IEEE International Symposium on Computer Architecture (**ISCA'20**). June 2020.
IEEE MICRO Top Picks Award
HiPEAC Paper Award
 26. Shang Wang, Yifan Bai, Gennady Pekhimenko.
Scaling Back-propagation by Parallel Scan Algorithm. Machine Learning and Systems Conference (**MLSys'20**). March 2020.
 25. Peter Mattson, Christine Cheng, Gregory Diamos, Cody Coleman, Paulius Micikevicius, David Patterson, Hanlin Tang, Gu-Yeon Wei, Peter Bailis, Victor Bittorf, David Brooks, Dehao Chen, Debo Dutta, Udit Gupta, Kim Hazelwood, Andy Hock, Xinyuan Huang, Daniel Kang, David Kanter, Naveen Kumar, Jeffery Liao, Deepak Narayanan, Tayo Oguntebi, Gennady Pekhimenko, Lillian Pentecost, Vijay Janapa Reddi, Taylor Robie, Tom St John, Carole-Jean Wu, Lingjie Xu, Cliff Young, Matei Zaharia.
MLPerf Training Benchmark. Machine Learning and Systems Conference (**MLSys'20**). March 2020.

24. Sihang Liu, Korakit Seemakhupt, Gennady Pekhimenko, Aasheesh Kolli, and Samira Khan.
Janus: Optimizing Memory and Storage Support for Non-Volatile Memory Systems. ACM/IEEE International Symposium on Computer Architecture (**ISCA'19**). June 2019.
IEEE MICRO Top Picks Honorable Mention
23. Hongyu Miao, Myeongjae Jeon, Gennady Pekhimenko, Kathryn S. McKinley, and Felix Xiaozhu Lin.
StreamBox-HBM: Stream Analytics on High Bandwidth Hybrid Memory. ACM International Conference on Architectural Support for Programming Languages and Operating Systems (**ASPLOS'19**). April 2019.
22. Anand Jayarajan, Jinliang Wei, Garth Gibson, Alexandra Fedorova, and Gennady Pekhimenko.
Priority-based Parameter Propagation for Distributed DNN Training. Machine Learning and Systems Conference (**MLSys'19**). April 2019.
21. Hongyu Zhu, Mohamed Akrouf, Bojian Zheng, Andrew Pelegris, Amar Phanishayee, Bianca Schroeder, and Gennady Pekhimenko.
Benchmarking and Analyzing Deep Neural Network Training. IEEE International Symposium on Workload Characterization (**IISWC'18**). October 2018.
20. Gennady Pekhimenko, Chuanxiong Guo, Myeongjae Jeon, Ryan Huang, and Lidong Zhou.
TerseCades: Efficient Data Compression in Stream Processing. USENIX Annual Technical Conference (**ATC'18**). July 2018.
19. Animesh Jain, Amar Phanishayee, Jason Mars, Lingjia Tang, and Gennady Pekhimenko.
Gist: Efficient Data Encoding for Deep Neural Network Training. International Symposium on Computer Architecture (**ISCA'18**). June 2018.
18. Nandita Vijaykumar, Abhilasha Jain, Diptesh Majumdar, Kevin Hsieh, Gennady Pekhimenko, Eiman Ebrahimi, Nastaran Hajinazaran, Phillip B. Gibbons, and Onur Mutlu .
A Case for Richer Cross-layer Abstractions: Bridging the Semantic Gap to Enhance Memory Optimization. International Symposium on Computer Architecture (**ISCA'18**). June 2018.
17. Hongyu Zhu, Bojian Zheng, Amar Phanishayee, Bianca Schroeder, and Gennady Pekhimenko.
DNN-Train: Benchmarking and Analyzing DNN Training. SysML Conference (**SysML'18**). February 2018.
16. Hongyu Miao, Heejin Park, Myeongjae Jeon, Gennady Pekhimenko, Kathryn S. McKinley, and Felix Xiaozhu Lin.
StreamBox: Modern Stream Processing on a Multicore Machine. USENIX Annual Technical Conference (**ATC'17**). July 2017.
15. Donghyuk Lee, Samira Khan, Lavanya Subramanian, Saugata Ghose, Rachata Ausavarungnirun, Gennady Pekhimenko, Vivek Seshadri, and Onur Mutlu.
Design-Induced Latency Variation in Modern DRAM Chips: Characterization, Analysis, and Latency Reduction Mechanisms. ACM SIGMETRICS / IFIP Performance (**SIGMETRICS'17**). June 2017.
14. Hasan Hassan, Nandita Vijaykumar, Samira Khan, Saugata Ghose, Kevin Chang, Gennady Pekhimenko, Donghyuk Lee, Oguz Ergin, and Onur Mutlu.
SoftMC: A Flexible and Practical Infrastructure for Enabling Experimental DRAM Studies. International Symposium on High-Performance Computer Architecture (**HPCA'17**). February 2017
13. Nandita Vijaykumar, Kevin Hsieh, Gennady Pekhimenko, Samira Khan, Ashish Shrestha, Saugata Ghose, Adwait Jog, Phillip B. Gibbons, Onur Mutlu.
Zorua: A Holistic Approach to Resource Virtualization in GPUs. International Symposium on Microarchitecture (**MICRO'16**). October 2016.
12. Kevin Chang, Abhijith Kashyap, Hasan Hassan, Saugata Ghose, Kevin Hsieh, Donghyuk

Lee, Tianshi Li, Gennady Pekhimenko, Samira Khan, Onur Mutlu.
Understanding Latency Variation in Modern DRAM Chips: Experimental Characterization, Analysis, and Optimization. ACM SIGMETRICS / IFIP Performance (**SIGMETRICS'16**). June 2016.

11. Gennady Pekhimenko, Evgeny Bolotin, Nandita Vijaykumar, Onur Mutlu, Todd C. Mowry, Stephen W. Keckler.
Toggle-Aware Bandwidth Compression for GPUs. International Symposium on High-Performance Computer Architecture (**HPCA'16**). March 2016.
10. Hasan Hassan, Gennady Pekhimenko, Nandita Vijaykumar, Vivek Seshadri, Donghyuk Lee, Oguz Ergin, Onur Mutlu.
ChargeCache: Reducing DRAM Latency by Exploiting Row Access Locality. International Symposium on High-Performance Computer Architecture (**HPCA'16**). March 2016.
9. Vivek Seshadri, Gennady Pekhimenko, Olatunji Ruwase, Onur Mutlu, Phillip B. Gibbons, Michael A. Kozuch, Todd C. Mowry, Trishul Chilimbi.
Page Overlays: An Enhanced Virtual Memory Framework to Enable Fine-grained Memory Management. International Symposium on Computer Architecture (**ISCA'15**). June 2015.
8. Nandita Vijaykumar, Gennady Pekhimenko, Adwait Jog, Abhishek Bhowmick, Rachata Ausavarungnirun, Onur Mutlu, Chita R. Das, Mahmut T. Kandemir, Todd C. Mowry.
A Case for Core-Assisted Bottleneck Acceleration in GPUs: Enabling Efficient Data Compression. International Symposium on Computer Architecture (**ISCA'15**). June 2015.
7. Gennady Pekhimenko, Dimitrios Lymberopoulos, Oriana Riva, Karin Strauss, Doug Burger.
PocketTrend: Architecting Search Engines for Trending Topics. International World Wide Web Conference (**WWW'15**). May 2015.
6. Gennady Pekhimenko, Tyler Hubery, Rui Cai, Onur Mutlu, Phillip B. Gibbons, Michael A. Kozuch, Todd C. Mowry.
Exploiting Compressed Block Size as an Indicator of Future Reuse. International Symposium on High-Performance Computer Architecture (**HPCA'15**). February 2015.
5. Donghyuk Lee, Yoongu Kim, Gennady Pekhimenko, Samira Khan, Vivek Seshadri, Kevin Chang, Onur Mutlu.
Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case. International Symposium on High-Performance Computer Architecture (**HPCA'15**). February 2015.
4. Bradley Thwaites, Gennady Pekhimenko, Amir Yazdanbakhsh, Girish Mururu, Jongse Park, Hadi Esmaeilzadeh, Onur Mutlu, Todd C. Mowry.
Rollback-Free Value Prediction with Approximate Loads. International Conference on Parallel Architectures and Compilation Techniques (**PACT'14, Short Paper**). August 2014.
3. Gennady Pekhimenko, Vivek Seshadri, Yoongu Kim, Hongyi Xin, Onur Mutlu, Phillip B. Gibbons, Michael A. Kozuch, Todd C. Mowry.
Linearly Compressed Pages: A Low-Complexity, Low-Latency Main Memory Compression Framework. International Symposium on Microarchitecture (**MICRO'13**). December 2013.
2. Vivek Seshadri, Yoongu Kim, Chris Fallin, Donghyuk Lee, Rachata Ausavarungnirun, Gennady Pekhimenko, Yixin Luo, Onur Mutlu, Phillip B. Gibbons, Michael A. Kozuch, Todd C. Mowry.
RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization. International Symposium on Microarchitecture (**MICRO'13**). December 2013.
1. Gennady Pekhimenko, Vivek Seshadri, Onur Mutlu, Phillip B. Gibbons, Michael A. Kozuch, Todd C. Mowry.
Base-Delta-Immediate Compression: Practical Data Compression Mechanism for On-Chip

Caches. International Conference on Parallel Architectures and Compilation Techniques (**PACT'12**). September 2012.

JOURNALS &
BOOK CHAPTERS

10. Anirudh Mohan Kaushik, Gennady Pekhimenko, Hiren Patel. *Gretch: A Hardware Prefetcher for Graph Analytics*. ACM Transactions on Architecture and Code Optimization (**TACO'21**). 2021.
9. Amir Yazdanbakhsh, Gennady Pekhimenko, Hadi Esmaeilzadeh, Onur Mutlu, Todd C. Mowry. *Towards Breaking the Memory Bandwidth Wall Using Approximate Value Prediction.. Approximate Circuits*. 2019.
8. Donghyuk Lee, Samira Manabi Khan, Lavanya Subramanian, Saugata Ghose, Rachata Ausavarungnirun, Gennady Pekhimenko, Vivek Seshadri, Onur Mutlu. *Design-Induced Latency Variation in Modern DRAM Chips: Characterization, Analysis, and Latency Reduction Mechanisms*. **POMACS: Proceedings of the ACM on Measurement and Analysis of Computing Systems**. 2017.
7. Hongyi Xin, Richard Zhu, Sunny Nahar, John Emmons, Gennady Pekhimenko, Carl Kingsford, Can Alkan, Onur Mutlu. *Optimal Seed Solver: Optimizing Seed Selection in Read Mapping*. **Oxford Bioinformatics**. 2016.
6. Amir Yazdanbakhsh, Gennady Pekhimenko, Bradley Thwaites, Girish Mururu, Jongse Park, Hadi Esmaeilzadeh, Onur Mutlu, Todd C. Mowry. *RFVP: Rollback-Free Value Prediction with Approximate Memory Loads*. ACM Transactions on Architecture and Code Optimization (**TACO'16**). 2016.
5. Donghyuk Lee, Saugata Ghose, Gennady Pekhimenko, Samira Khan, Onur Mutlu. *Simultaneous Multi Layer Access: A High Bandwidth and Low Cost 3D-Stacked Memory Interface*. ACM Transactions on Architecture and Code Optimization (**TACO'16**). 2015.
4. Amir Yazdanbakhsh, Gennady Pekhimenko, Bradley Thwaites, Girish Mururu, Jongse Park, Hadi Esmaeilzadeh, Onur Mutlu, Todd C. Mowry. *Mitigating the Bandwidth Bottleneck with Approximate Load Value Prediction*. **IEEE Design & Test**. 2016.
3. Gennady Pekhimenko, Evgeny Bolotin, Mike O'Connor, Onur Mutlu, Todd C. Mowry, Stephen W. Keckler. *Toggle-Aware Compression for GPUs*. IEEE Computer Architecture Letters (**CAL'15**). May 2015.
2. Hongyi Xin, John Greth, John Emmons, Gennady Pekhimenko, Carl Kingsford, Can Alkan, Onur Mutlu. *Shifted Hamming Distance: A Fast and Accurate SIMD-Friendly Filter for Local Alignment in Read Mapping*. **Oxford Bioinformatics**. January 2015.
1. Gennady Pekhimenko, Angela Demke Brown. *Software Automatic Tuning: From Concepts to State-of-the-Art Results, Chapter 19*. **Springer**. September 2010.

OTHER
PEER-REVIEWED
PUBLICATIONS

4. Bojian Zheng and Gennady Pekhimenko. *EcoRNN: Efficient Computing of LSTM RNN on GPUs*. Student Research Competition at IEEE/ACM International Symposium on Microarchitecture (**SRC@MICRO'18**). October 2018.
3. Gennady Pekhimenko, Evgeny Bolotin, Mike O'Connor, Onur Mutlu, Todd C. Mowry, Stephen W. Keckler. *Energy-Efficient Data Compression for GPU Memory Systems*. Student Research Competition at International Conference on Architectural Support for Programming Languages and Operating Systems (**SRC@ASPLOS'15**). March 2015.

2. Gennady Pekhimenko, Todd C. Mowry, Onur Mutlu.
Linearly Compressed Pages: A Main Memory Compression Framework with Low Complexity and Low Latency. Student Research Competition at International Conference on Parallel Architectures and Compilation Techniques (**SRC@PACT'12**). September 2012.

1. Gennady Pekhimenko, Angela Demke Brown.
Efficient Program Compilation through Machine Learning Techniques. International Workshop on Automatic Performance Tuning (**iWAPT'09**). October 2009

PATENTS, THESES

6. Amar Phanishayee, Gennady Pekhimenko, Animesh Jain.
Efficient data encoding for deep neural network training. Patent No. 20190347549. November 2019.
5. Gennady Pekhimenko.
Practical Data Compression for Modern Memory Hierarchies. PhD Thesis, Carnegie Mellon University. July 2016.
4. Dimitrios Lymberopoulos, Oriana Riva, Karin Strauss, Doug Burger, Gennady Pekhimenko.
Trend Response Management. Patent No. 20150227517. August 2015.
3. Yaoqing Gao, Tong Chen, Zehra Sura, Gennady Pekhimenko, Kevin O'Brien, Khaled Mohammed, Roch Archambault, Raul Silvera.
Managing Speculative Assist Threads. Patent No. 20110093838. October 2010.
2. Gennady Pekhimenko.
Machine Learning Algorithms for Choosing Compiler Heuristics. MS Thesis, University of Toronto. January 2008.
1. Gennady Pekhimenko.
Performance Analysis of MPI-Programs. Diploma Thesis, Moscow State University, Russia. May 2004.

GRANTS

- ◇ NSERC USRA Award “Machine Learning Compilers”, USRA: Benjamin Chislett, **\$6,000 total**. 2021–2021
- ◇ DCS Award “Methodology for Developing and Evaluating Machine Learning Chips”, DCS: Chenhao Jiang, **\$6,000 total**. 2021–2021
- ◇ Mitacs, Accelerate, “Adversarial Robustness of Deep Learning Algorithms on Next-Gen AI Accelerators”, **\$30,000 total**. 2021–2022
- ◇ Canada Foundation for Innovation (CFI), John Evans Leaders Fund program, Co-PI, “Computer systems support for machine learning and artificial intelligence”, Total: **CAD\$276,000**, My share: **CAD\$138,000**. 2021–2024
- ◇ Amazon, AWS Machine Learning Research Award, “Efficient DNN Training at Scale: from Algorithms to Hardware”, **USD\$40,000 cash and USD\$80,000 in AWS cloud credits**. 2020–2021
- ◇ Facebook, Faculty Research Award (AI Systems HW/SW Co-Design), “Efficient DNN Training at Scale: from Algorithms to Hardware”, **USD\$49,500**. 2020–2021
- ◇ Facebook, Facebook/University of Toronto unrestricted gift, Co-PI, “Efficient ML Everywhere: From the Edge to the Data Center, From SW to HW”, Total: **USD\$150,000**, My share: **USD\$50,000**. 2020–2021
- ◇ Mitacs, Accelerate, “Explore efficiently automated parallel hyperparameter search for optimizing machine learning models over large scale cloud cluster”, **\$30,000 total**. 2020–2021
- ◇ NSERC UTEA “Fair and Efficient Scheduling of Machine Learning Workloads in High Performance Computer Clusters”, UTEA: Yu Bo Gao, **\$4,875 total**. 2020–2020
- ◇ ESROP - U of T “Efficient Streaming Engines for Time Series Data”, ESROP: Kimberly Hau, **\$3,000 total**. 2020–2020
- ◇ NSERC CRD, “Efficient Distributed DNN Training and Inference”, **\$273,000 total**.

2020–2023

- ◇ NSERC CRD, “Efficient Compiler-Driven Pointer Compression”, **\$180,000 total**. 2020–2023
- ◇ CIFAR, AI Chair, “Systems for Machine Learning”, **\$1,000,000**. 2019–2024
- ◇ NSERC, Strategic Networks Grant, Co-PI, “COHESA Network”, Total: **\$1,000,000 per year**, My share: **\$35,000 per year**. 2019–2021
- ◇ Mitacs, Accelerate, “Next Generation AI Accelerator Algorithm Hardware Co-Optimization”, **\$30,000 total**. 2019–2020
- ◇ Huawei, Research Grant, “Efficient Data Compression/Deduplication for Persistent Memory and DRAM”, **\$428,400 total**. 2019–2022
- ◇ IBM Canada, CAS Program, “Efficient Compiler-Driven Pointer Compression”, Award #1112, **\$90,000 total**. 2019–2022
- ◇ Huawei, Research Grant, “Compiler Infrastructure for Optimizing DNN Workloads”, **\$289,300 total**. 2019–2022
- ◇ NSERC Discovery Grant (Increase) “Exploiting Hardware Heterogeneity for Efficient Execution of Emerging Applications”, **\$12,500 total**. 2018–2023
- ◇ NSERC CRD, “Efficient Memory Footprint Reduction for Java Performance”, **\$204,000 total**. 2019–2022
- ◇ Huawei, Research Grant, “Efficient Distributed DNN Training”, **\$199,546 total**. 2018–2021
- ◇ NSERC UTEA “Parallelism and Hardware Heterogeneity Support in Modern Compilers”, UTEA: Qiongsi Wu, **\$4,875 total**. 2018–2018
- ◇ Connaught New Researcher Award “Exploiting Hardware Heterogeneity for Efficient Execution of Emerging Applications”, Connaught Fund, **\$10,000 total**. 2018–2019
- ◇ NSERC Discovery Accelerator Supplement Grant “Exploiting Hardware Heterogeneity for Efficient Execution of Emerging Applications”, NSERC (522575), **\$120,000 total**. 2018–2021
- ◇ NSERC Discovery Grant “Exploiting Hardware Heterogeneity for Efficient Execution of Emerging Applications”, NSERC (RGPIN-2018-06514), **\$140,000 total**. 2018–2023
- ◇ IBM Canada, CAS Program, “Efficient Memory Footprint Reduction for Java Performance”, Award #1063, **\$102,000 total**. 2018–2021
- ◇ Huawei, HiRP Open Program, “Hardware/Software Optimization and Compiler Support for Heterogeneous Systems”, **\$87,044 total**. 2017–2019
- ◇ Canada Foundation for Innovation (CFI), John Evans Leaders Fund program “Heterogeneous Systems Laboratory”, CFI (Award #36585), **\$240K total**. 2017–2020

STUDENTS

Student Award Highlights: Geoffrey Yu (NSERC CGS-D, NSERC CGS-M, Snap Research Scholarship, Vector Institute Scholarship in Artificial Intelligence, Queen Elizabeth II Graduate Scholarship), Hanjie Qiu (OGS, Vector Institute Scholarship in Artificial Intelligence), Qiongsi Wu (NSERC CGS-M, Vector Institute Scholarship in Artificial Intelligence), Serina Tan (Vector Institute Scholarship in Artificial Intelligence), Bojian Zheng (Third place in MICRO 2018 ACM SRC), James Gleeson (Bell Scholarship (twice)), Qidong Su (Wolfand Scholarship), Kevin Song (OGS, Vector Institute Scholarship in Artificial Intelligence).

Student Supervising Summary: Currently, supervise 9 PhD students (1 co-advised), and 7 Masters students. Graduated: 10 Masters students.

Current:

- ◇ Hongyu Zhu, PhD Student. DNN Profiling and Analysis.
- ◇ Bojian Zheng, PhD student. Hardware acceleration for LSTM-based RNNs.
- ◇ Alexandra Tsvetkova, PhD Student. Software support for GPU virtualization.
- ◇ James Gleeson, PhD Student (co-advised with Eyal de Lara). Optimizing reinforcement learning training.
- ◇ Anand Jayarajan, PhD Student. Efficient stream processing engine.

- ◇ Mustafa Quraish, PhD student.
- ◇ Shang (Sam) Wang, PhD student. Horizontal fusion for efficient DNN training.
- ◇ Jiacheng Yang, PhD student. DNN Training at the edge.
- ◇ Qidong Su, PhD student.
- ◇ Pavel Golikov, MSc. student. DNNs performance modeling on modern GPUs.
- ◇ Yaoyao Ding, MSc. student.
- ◇ Kevin Song, MSc. student.
- ◇ Daniel Snider, MSc. student.
- ◇ Jasper Zhu, MSc. student.
- ◇ Peiming Yang, MSc. student
- ◇ Xin Li, MSc. student

- ◇ Yu Bo Gao, BSc. student.
- ◇ Chenhao Jiang, BSc. student.
- ◇ Zhanda Zhu, BSc. student from SJTU.
- ◇ Qingyuan Qiu, BSc. student.
- ◇ Wei Zhao, BSc. student.
- ◇ Murali Andoor, BSc. student.
- ◇ Yudi Sun, BSc. student.
- ◇ Benjamin Chislett, BSc. student.
- ◇ Yvonne Yang, BSc. student
- ◇ Maryam Gohargani, BSc. student

Graduated:

- ◇ Qionsgi Wu, MSc. (2021). Compiler support for multi-threading with OpenMP. First position: IBM.
- ◇ Hanjie Qiu, MSc. (2021). Pavise: Integrating Fault Tolerance Support for Persistent Memory Applications. First position: industry.
- ◇ Xiaodan (Serina) Tan, MSc. (2021). GPUPool: A Holistic Approach to Fine-Grained GPU Sharing in the Cloud. First position: Amazon AWS.
- ◇ Jiahuang (Jacob) Lin, MScAC. (2021). Speech recognition using DeepSpeech2 model.
- ◇ Geoffrey Yu, MSc. (2020). Habitat: Prediction-guided Hardware Selection for Deep Neural Network Training. First position: PhD student at MIT EECS.
- ◇ Shang (Sam) Wang, MSc. (2020). Back-propagation by Parallel Scan Algorithm. First position: Nvidia.
- ◇ Yingying Fu, MScAC (2020). Next Generation AI Accelerator Algorithm Hardware Co-Optimization. First position: Untether AI, Toronto, ON.
- ◇ Izaak Niksan, BSc. (2020). Memory profiler for DNN training.
- ◇ Pavel Klishin, MSc. (2019). DNN training acceleration with FPGAs. First position: Industry, Moscow, Russia
- ◇ Andrew Pelegris, MSc. (2019) . Binarized DNNs acceleration. First position: Stealth-mode startup.
- ◇ Mohamed Akrouf, MScAC (2019). Reinforcement learning profiling and training. First position: Research Scientist at Triage, Toronto, ON.
- ◇ Ming (Michael) Yang, BSc. (2019). New simulator infrastructure for GPUs. First position: Engineer at Cerebras, Bay Area, CA.
- ◇ Yifan Bai, BSc. (2019). Jacobian-based approach for DNN training. First position: Graduate student at UC Berkeley, CA.
- ◇ Kuei-Fang (Albert) Hsueh, BSc. (2019). Machine translation using Transformer model for inference. First position: Graduate student at UofT, Toronto, ON
- ◇ Akshay Nair, BSc. (2018). Simulation infrastructure for GPUs. First position: Software Engineer at Google, Mountain View.

MENTORING CMU (PhD, Masters and undergraduate):

- ◇ Amir Yazdanbakhsh, PhD Student. Research Project: Rollback-Free Value Prediction

with Approximate Loads.

- ◇ Hasan Hassan, Masters student. Research Project: Reducing DRAM Latency by Exploiting Row Access Locality.
- ◇ Mahmoud Khairy, Masters student. Research Project: Efficient DRAM Refresh for GPUs.
- ◇ Arthur Perais, PhD Student. Research Project: Synergy Analysis Between Value Prediction and Data Compression.
- ◇ Hongyi Xin, PhD Student. Research Project: Shifted Hamming Distance: A Fast and Accurate SIMD-Friendly Filter for Local Alignment in Read Mapping.
- ◇ Nandita Vijaykumar, PhD Student. Research Project: Core-Assisted Bottleneck Acceleration.
- ◇ Abhishek Bhowmick, Undergraduate student (currently Masters student at CMU). Research Project: GPU Main Memory Compression and Prefetching.
- ◇ Tyler Huberty and Rui Cai, Undergraduate students (currently at Apple and Microsoft). Research Project: CARP: Compression-Aware Replacement Policies.
- ◇ Jason Lin and Brian Osbun, Undergraduate students (currently at Microsoft and CMU). Research Project: Bandwidth-Optimized Prefetching.
- ◇ Martyn Romanko and Lei Fan, Masters students (currently at Intel). Research Project: Implementation and Energy Analysis of Base-Delta-Immediate Compression.

WORK
EXPERIENCE

- ◇ Faculty Member at **Vector Institute**, Sep 2019 – Present
- ◇ Assistant Professor at the **University of Toronto**, CS Department, July 2017 – Present
- ◇ Assistant Professor at the **University of Toronto**, ECE Department (by courtesy), January 2018 – Present
- ◇ Researcher at **Microsoft Research** with Systems Research Group, July 2016 – Aug 2017
- ◇ Graduate Student Researcher at **Carnegie Mellon University** with Prof. Todd C. Mowry and Prof. Onur Mutlu, Sep 2010 – August 2016
- ◇ Research Consultant at **Microsoft** with Dr. Marc Tremblay, Feb 2015 – Jul 2015
- ◇ Research Intern at **NVIDIA Research** with Dr. Stephen Keckler and Dr. Evgeny Bolotin, Summer 2014
- ◇ Research Intern at **Microsoft Research** with Dr. Karin Strauss, Dr. Dimitrios Lybmeropoulos, Dr. Oriana Riva, Summer 2013
- ◇ Research Intern at **Microsoft Research** with Dr. Ella Bounimova, Dr. Patrice Godefroid, and Dr. David Molnar, Summer 2012
- ◇ Compiler Engineer/Researcher (Full-time) at **IBM** with Raul Silvera and Yaoging Gao, May 2007 – Jun 2010
- ◇ Graduate Student Researcher at the **University of Toronto** with Prof. Angela Demke Brown, Sep 2006 – Jan 2008
- ◇ Compiler Engineer (Full-time) at **Elbrus, Moscow, Russia** with Dr. Vladimir Volkonskii, May 2004 – Aug 2006
- ◇ System Programmer at **Intel-MSU Lab, Moscow, Russia** with Prof. Viktor Krukov, May 2003 – Jun 2004

TEACHING
EXPERIENCE

- ◇ **Instructor** at the University of Toronto, Fall 2020
CSC B58H: Computer Organization, Undergraduate
- ◇ **Instructor** at the University of Toronto, Winter 2021, 2020, 2019, 2018
CSC D70H: Compiler Optimization, Undergraduate
- ◇ **Instructor** at the University of Toronto, Fall 2021, 2020, 2019, 2018, 2017
CSC 2224H: Parallel Computer Architecture and Programming, Graduate
- ◇ **Teaching Assistant** at Carnegie Mellon University, Spring 2012
Optimizing Compilers, Graduate
- ◇ **Teaching Assistant** at Carnegie Mellon University, Fall 2011
Introduction to Computer Systems, Undergraduate

	◇ Teaching Assistant at the University of Toronto, Operating Systems, Undergraduate	Fall 2007, Spring 2008
	◇ Teaching Assistant at the University of Toronto, Computer Programming, Undergraduate	Spring 2007
	◇ Teaching Assistant at the University of Toronto, Software Engineering, Undergraduate	Fall 2006
INVITED TALKS	67. <i>Efficient DNN Training at Scale</i> SAFARI @ETH Zurich, online	August 2021
	66. <i>LifeStream: A High-Performance Stream Processing Engine for Periodic Streams</i> Microsoft Research, online	May 2021
	65. <i>Apple ODML Workshop Keynote Invited Talk</i> Apple ODML Workshop, online	April 2021
	64. <i>ASPLOS MLBench Workshop</i> ASPLOS MLBench'21, online	April 2021
	63. <i>MLSys MLBench Tutorial</i> MLSys MLBench'21, online	April 2021
	62. <i>HPCA MLBench Tutorial</i> HPCA MLBench'21, online	February 2021
	61. <i>Efficient DNN Training at Scale: from Algorithms to Hardware</i> Facebook, Facebook Faculty Summit, online	October 2020
	60. <i>Keynote talk on Efficient DNN Training at Scale</i> Vector Institute NLP Symposium, online	September 2020
	59. <i>ISPASS Tutorial on ML Benchmarking</i> ISPASS ML Performance'20, online	August 2020
	58. <i>VCEW Invited Talk: ML Benchmarking</i> VCEW'20, online	June 2020
	57. <i>ISCA Mini-Panel: Accelerators</i> ISCA'20, online	June 2020
	56. <i>Efficient DNN Training at Scale: from Algorithms to Hardware</i> Microsoft, Microsoft Research Seminar, online	May 2020
	55. <i>Efficient DNN Training at Scale: from Algorithms to Hardware</i> Facebook, SysML Seminar, online	May 2020
	54. <i>Holistic Approach to DNN Training Efficiency: Analysis and Optimizations</i> Fields Institute, Toronto, ON	January 2020
	53. <i>Holistic Approach to DNN Training Efficiency: Analysis and Optimizations</i> Uber ATG, Toronto, ON	November 2019
	52. <i>Holistic Approach to DNN Training Efficiency: Analysis and Optimizations</i> Yandex, Moscow, Russia	August 2019
	51. <i>ML Performance: Benchmarking Deep Learning Systems</i> ISCA'19 Tutorial, Phoenix, Ar.	June 2019
	50. <i>ML Performance: Benchmarking Deep Learning Systems</i> ASPLOS'19 Tutorial, Providence, RI.	April 2019
	49. <i>Holistic Approach to DNN Training Efficiency: Analysis and Optimizations</i> Google Platform Team, Sunnyvale, CA.	April 2019
	48. <i>Holistic Approach to DNN Training Efficiency: Analysis and Optimizations</i> FastPath'19 Workshop Keynote, Madison, WI.	March 2019
	47. <i>Holistic Approach to DNN Training Efficiency: Analysis and Optimizations</i> Apple, Cupertino, CA.	December 2018
	46. <i>Holistic Approach to DNN Training Efficiency: Analysis and Optimizations</i> Facebook, Menlo Park, CA.	December 2018
	45. <i>Holistic Approach to DNN Training Efficiency: Analysis and Optimizations</i> Google, Mountain View, CA.	December 2018

44. *Algorithms vs. Architectures: Rivals or Partners in Pushing AI Boundaries?*
Huawei AI Workshop, Shanghai, China. October 2018
43. *TerseCades: Efficient Data Compression in Stream Processing*
USENIX ATC'18, Boston, MA. July 2018
42. *Benchmarking and Analyzing DNN Training*
Vector Institute, Toronto, ON. May 2018
41. *Benchmarking and Analyzing DNN Training*
SysML'18, Stanford, CA. Feb 2018
40. *Practical Data Compression for Memory Hierarchy and DNNs*
Yandex, Moscow, Russia. Nov 2017
39. *A Case for Toggle-Aware Compression for GPU Systems*
HPCA-22, Barcelona, Spain. Mar 2016
38. *RFVP: Rollback-Free Value Prediction with Safe-to-Approximate Loads*
HiPEAC16, Prague, Czech Republic. Jan 2016
37. *Linearly Compressed Pages*
University of Texas at Austin, Austin, TX. Nov 2015
36. *Exploiting Compressed Block Size as an Indicator of Future Reuse*
PDL Retreat, Bedford, PA. Oct 2015
35. *Linearly Compressed Pages*
University of Illinois at Urbana-Champaign, Urbana, IL. Oct 2015
34. *Base-Delta-Immediate Compression*
University of Alberta, Edmonton, Canada. Sep 2015
33. *PocketTrend: Timely Identification and Delivery of Trending Search Content to Mobile Users*
WWW-24, Florence, Italy. May 2015
32. *Exploiting Compressed Block Size as an Indicator of Future Reuse*
MIT, Boston, MA. May 2015
31. *Energy-Efficient Data Compression for Modern Memory Systems*
QInF Finals, San Diego, CA. Mar 2015
30. *Energy-Efficient Data Compression for GPU Memory Systems*
SRC@ASPLOS'15, Istanbul, Turkey. **First Place in ACM SRC Competition** Mar 2015
29. *Exploiting Compressed Block Size as an Indicator of Future Reuse*
Intel Atom Group, Hillsboro, OR. Feb 2015
28. *Exploiting Compressed Block Size as an Indicator of Future Reuse*
HPCA-21, Bay Area, CA. Feb 2015
27. *Energy-Efficient Data Compression*
Qualcomm, San Diego, CA. Sep 2014
26. *Energy-Efficient Data Compression For GPU Memory Systems*
NVIDIA Research, Santa Clara, CA. Sep 2014
25. *Linearly Compressed Pages*
Intel Labs, Santa Clara, CA. Sep 2014
24. *Energy-Efficient Data Compression*
UC Berkeley, ASPIRE Lab, Berkeley, CA. Sep 2014
23. *Linearly Compressed Pages*
Huawei R&D, Santa Clara, CA. Aug 2014
22. *Energy-Efficient Data Compression*
NVIDIA Research, Santa Clara, CA. Jul 2014
21. *Guest Lecture on Cache Compression*
18447: Introduction to Computer Architecture, Pittsburgh, PA. Apr 2014
20. *Linearly Compressed Pages*
CMU Cloud Workshop, Pittsburgh, PA. Apr 2014
19. *Linearly Compressed Pages*
Samsung Research, San Jose, CA. Dec 2013
18. *Linearly Compressed Pages*

	Oracle Labs, Belmont, CA.	Dec 2013
17.	<i>Linearly Compressed Pages: A Low-Complexity, Low-Latency Main Memory Compression Framework</i> MICRO-46, Davis, CA.	Dec 2013
16.	<i>Linearly Compressed Pages</i> Stanford Cloud Workshop, Mountain View, CA.	Dec 2013
15.	<i>Linearly Compressed Pages</i> NVIDIA Research, Santa Clara, CA.	Dec 2013
14.	<i>Main Memory Compression and Low-Cost Compression Algorithms</i> PDL Retreat, Bedford, PA.	Oct 2013
13.	<i>In-Memory Optimizations: Efficient Compression and Data Movement</i> Heidelberg Laureate Forum, Heidelberg, Germany.	Sep 2013
12.	<i>PocketTrend: Efficient Trend Detection for Mobile Devices</i> Microsoft Research, Redmond, WA.	Aug 2013
11.	<i>Base-Delta-Immediate Compression</i> Microsoft Research, Redmond, WA.	Jul 2013
10.	<i>Base-Delta-Immediate Compression</i> University of Toronto, Ontario, Canada.	Mar 2013
9.	<i>Heterogeneous Block Architectures</i> Qualcomm, San Diego, CA.	Mar 2013
8.	<i>Linearly Compressed Pages</i> Intel, Hillsboro, OR.	Feb 2013
7.	<i>Base-Delta-Immediate Compression</i> Intel Labs, Hillsboro, OR.	Feb 2013
6.	<i>Guest Lecture on Caching in Multi-Core Systems</i> 18742: Parallel Computer Architecture, Pittsburgh, PA.	Oct 2012
5.	<i>Base-Delta-Immediate Compression: Practical Data Compression Mechanism for On-Chip Caches</i> PACT, Minneapolis, MN.	Sep 2012
4.	<i>Linearly Compressed Pages: A Main Memory Compression Framework with Low Complexity and Low Latency</i> SRC@PACT, Minneapolis, MN. Second Place in ACM SRC Competition	Sep 2012
3.	<i>Guest Lecture on Dynamic Compilation</i> 15745: Optimizing Compilers, Pittsburgh, PA.	Feb 2012
2.	<i>Assist Threads for Data Prefetching in IBM XL Compilers</i> CASCON, Toronto, ON.	Nov 2009
1.	<i>Efficient Program Compilation through Machine Learning Techniques</i> International Workshop on Automatic Performance Tuning, Tokyo, Japan.	Oct 2009

SERVICE

MLCommons/MLPerf Research Co-Chair 2019–2021

Program and Organization Committees

◇ Heavy PC Member, EuroSys 2022	2021–2022
◇ PC Member, MLSys 2022	2021–2022
◇ PC Member, ASPLOS 2022	2021–2022
◇ PC Member, MICRO 2021	2021–2021
◇ Co-Chair, Artifact Evaluation Committee at MICRO 2021	2021–2021
◇ PC Member, OSDI 2021	2020–2021
◇ ERC (External Review Committee) Member, ISCA 2021	2020–2021
◇ Chair, Artifact Evaluation Committee at ASPLOS 2021	2020–2021
◇ PC Member, MLSys 2021	2020–2021
◇ PC Member, HPCA 2021	2020–2021
◇ PC Member, MICRO 2020	2020
◇ ERC (External Review Committee) Member, ISCA 2020	2019–2020

	◇ PC Member , MICRO TopPicks 2020	2019–2020
	◇ PC Member , MLSys 2020	2019–2020
	◇ PC Member , EuroSys 2020	2019–2020
	◇ Co-Chair , Artifact Evaluation at MLSys 2020	2019–2020
	◇ Publicity Co-Chair , HPCA 2020	2019–2020
	◇ PC Member , CGO 2020	2019–2020
	◇ PC Member , MICRO 2019	2019
	◇ PC Member , ICS 2019	2018–2019
	◇ Tutorial Organizer , MLPerfBench at ISCA 2019	2019
	◇ PC Member , ISCA 2019	2018–2019
	◇ PC Member , MLSys 2019	2018–2019
	◇ Co-Chair , Artifact Evaluation at SysML 2019	2018–2019
	◇ ERC (External Review Committee) Member , ASPLOS 2019	2018–2019
	◇ ERC (External Review Committee) Member , HPCA 2019	2018–2019
	◇ Program Co-Chair , Compiler-Driven Performance Workshop	2018
	◇ PC Member , MICRO 2018	2018
	◇ PC Member , ICS 2018	2017–2018
	◇ Publicity Co-Chair , ASPLOS 2018	2017–2018
	◇ ERC (External Review Committee) Member , MICRO 2017	2017
	◇ ERC (External Review Committee) Member , ISCA 2017	2016–2017
	◇ Web Chair , ISCA 2017	2016–2017
	◇ PC Member , ICWE 2017	2016–2017
	◇ ERC (External Review Committee) Member , ISCA 2016	2015–2016
	◇ PC Member , WWW 2016	2015–2016
	◇ Web Chair , ASPLOS 2016	2015–2016
	◇ Publicity Chair , HiPEAC 2015	2015
	◇ Information Director , Transactions on Computer Systems (TOCS)	2013–2017
PROFESSIONAL MEMBERSHIPS	◇ IEEE Computer Society	2014–present
	◇ Association of Computing Machinery (ACM)	2012–present
	◇ ACM SIGARCH	2012–present
CITIZENSHIP	◇ Russian Citizenship	
	◇ Canadian Permanent Resident	