

Contact Information	<p>School of Computing KAIST 291 Daehak-ro, Yuseong-gu Daejeon, South Korea, 34141</p> <p><i>E-mail:</i> jspark@casys.kaist.ac.kr <i>URL:</i> https://jongse-park.github.io</p>
Research Interests	Computer Architecture, Computer System, HW/SW Co-Design, Generative AI Serving Systems, On-Device AI Systems, Processing-in-Memory
Employment	<p>Visiting Associate Professor. Stanford University Jan. 2025–date Associate Professor. KAIST Mar. 2024–date Assistant Professor. KAIST Dec. 2019–Feb. 2024 System Architect. Bigstream Solutions Inc. Jun. 2018–Nov. 2019</p>
Education	<p>Ph.D. in Computer Science. Georgia Institute of Technology Aug. 2013–Aug. 2018 <ul style="list-style-type: none"> • Advisor: Prof. Hadi Esmaeilzadeh • Dissertation: <i>Breaking the Abstractions for Productivity and Performance in the Era of Specialization</i> </p> <p>M.S. in Computer Science. KAIST Feb. 2012 <ul style="list-style-type: none"> • Advisor: Prof. Seungryoul Maeng • Thesis: <i>Dynamic Resource Reconfiguration on the Cloud for Improving Data Locality</i> </p> <p>B.E. in Computer Science and Engineering. Sogang University Feb. 2010 <ul style="list-style-type: none"> • Graduated with Honors </p>
Honors and Awards	<p>Samsung Humantech Paper Award. 2025 Gold Prize (1st place in the Computer Science and Engineering track)</p> <p>Best Paper Award & Distinguished Artifact Award. IISWC. 2024 “LLMServingSim: A HW/SW Co-Simulation Infrastructure for LLM Inference Serving at Scale”</p> <p>Distinguished Artifact Award. ISCA. 2024 “DACAPO: Accelerating Continuous Learning in Autonomous Systems for Video Analytics”</p> <p>ISCA 25-Year Retrospective 1996-2020 Inclusion 2023 “Bit Fusion: Bit-Level Dynamically Composable Architecture for Accelerating Deep Neural Networks”</p> <p>ISCA 25-Year Retrospective 1996-2020 Inclusion 2023 “General-Purpose Code Acceleration with Limited-Precision Analog Computation”</p> <p>Distinguished Paper Award. HPCA. 2016 “TABLA: A Unified Template-Based Framework for Accelerating Statistical Machine Learning”</p> <p>Honorable Mention in IEEE Micro Top Picks from 2014 Computer Architecture Conferences. 2015 “General-Purpose Code Acceleration with Limited-Precision Analog Computation”</p>
Refereed Conference Papers	<ol style="list-style-type: none"> 1. M. Kim, S. Hong, R. Ko, S. Choi, H. Lee, J. Kim, J-Y Kim, J. Park, “Oaken: Fast and Efficient LLM Serving with Online-Offline Hybrid KV Cache Quantization,” in <i>International Symposium on Computer Architecture (ISCA)</i>, June 2025. (To Appear) 2. Y. Kim, I. Kim, K. Choi, J. Ahn, J. Park, J. Huh, “Interference-Aware DNN Serving on Heterogeneous Processors in Edge Systems,” in <i>IEEE International Conference on Computer Design (ICCD)</i>, November 2024. 3. J. Cho, M. Kim, H. Choi, G. Heo, J. Park, “LLMServingSim: A HW/SW Co-Simulation Infrastructure for LLM Inference Serving at Scale,” in <i>International Symposium on Workload Characterization (IISWC)</i>, September 2024.

4. M. Kim, J. Hwang, G. Heo, S. Cho, D. Mahajan, **J. Park**, "Accelerating String-key Learned Index Structures via Memoization-based Incremental Training," in *International Conference on Very Large Data Bases (VLDB)*, August 2024.
5. Y. Kim, C. Oh, J. Hwang, W. Kim, S. Oh, Y. Lee, H. Sharma, A. Yazdanbakhsh, **J. Park**, "DaCapo: Accelerating Continuous Learning in Autonomous Systems for Video Analytics," in *International Symposium on Computer Architecture (ISCA)*, June 2024.
6. G. Heo, S. Lee, J. Cho, H. Choi, S. Lee, H. Ham, G. Kim, D. Mahajan, **J. Park**, "NeuPIMs: NPU-PIM Heterogeneous Acceleration for Batched LLM Inferencing," in *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 2024.
7. S. Ghodrati, S. Kinzer, H. Xu, R. Mahapatra, Y. Kim, B. H. Ahn, D. K. Wang, L. Karthikeyan, A. Yazdanbakhsh, **J. Park**, N. S. Kim, H. Esmaeilzadeh, "Tandem Processor: Grappling with Emerging Operators in Neural Networks," in *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 2024.
8. Sunho Lee, Seonjin Na, Jungwoo Kim, Jongse Park, and Jaehyuk Huh, "Tunable Memory Protection for Secure Neural Processing Units," in *The 40th IEEE International Conference on Computer Design (ICCD)*, October 2022.
9. Bokyeong Kim, Soojin Hwang, Sanghoon Cha, Chang Hyun Park, Jongse Park, and Jaehyuk Huh, "Supporting Dynamic Translation Granularity for Hybrid Memory Systems," in *The 40th IEEE International Conference on Computer Design (ICCD)*, October 2022.
10. Joon Kyung Kim, Byung Hoon Ahn, Sean Kinzer, Soroush Ghodrati, Rohan Mahapatra, Brahmendra Yatham, Dohee Kim, Parisa Sarikhani, Babak Mahmoudi, Divya Mahajan, Jongse Park, Hadi Esmaeilzadeh, "Yin-Yang: Programming Abstraction for Cross-Domain Multi-Acceleration," in *IEEE Micro, special issue on Compiling for Accelerators*, 2022.
11. Jinwoo Hwang, Minsu Kim, Daeun Kim, Seungho Nam, Yoonsung Kim, Dohee Kim, Hardik Sharma, Jongse Park, "CoVA: Exploiting Compressed-Domain Analysis to Accelerate Video Analytics," in *USENIX Annual Technical Conference (ATC)*, July 2022.
12. Seungbeom Choi, Sunho Lee, Yeonjae Kim, Jongse Park, Youngjin Kwon, and Jaehyuk Huh, "Serving Heterogeneous Machine Learning Models on Multi-GPU Servers with Spatio-Temporal Sharing," in *USENIX Annual Technical Conference (ATC)*, July 2022.
13. S. Lee, J. Kim, S. Na, **J. Park**, and J. Huh, "TNPU: Supporting Trusted Execution with Tree-less Integrity Protection for Neural Processing Unit," in *The 27th IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, February 2022. [To appear]
14. S. Na, S. Lee, Y. Kim, **J. Park**, and J. Huh, "Common Counters: Compressed Encryption Counters for Secure GPU Memory," in *The 27th IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, February 2021.
15. S. Ghodrati, H. Sharma, S. Kinzer, A. Yazdanbakhsh, **J. Park**, N. Kim, D. Burger, and H. Esmaeilzadeh, "Mixed-Signal Charge-Domain Acceleration of Deep Neural Networks through Interleaved Bit-Partitioned Arithmetic," in *The 29th International Conference on Parallel Architectures and Compilation Techniques (PACT)*, October 2020.
16. Y. Li, **J. Park**, M. Alian, Y. Yuan, Q. Zheng, P. Pan, R. Wang, A. Schwing, H. Esmaeilzadeh, N. Kim, "A Network-Centric Hardware/Algorithm Co-Design to Accelerate Distributed Training of Deep Neural Networks," *The 50th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, October 2018.
17. H. Sharma, **J. Park**, B. Samynathan, B. Robatmili, S. Mirkhani, H. Esmaeilzadeh, "From Tensors to FPGAs: Accelerating Deep Learning," *A Symposium on High Performance Chips (Hot Chips)*, August 2018.
18. H. Sharma, **J. Park**, N. Suda, L. Lai, B. Chau, J. Kim, V. Chandra, H. Esmaeilzadeh, "Bit Fusion: Bit-Level Dynamically Composable Architecture for Accelerating Deep Neural Networks," *International Symposium on Computer Architecture (ISCA)*, June 2018.

19. **J. Park**, H. Sharma, D. Mahajan, J. Kim, P. Olds, H. Esmaeilzadeh, "Scale-Out Acceleration for Machine Learning," in *The 50th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, October 2017.
20. **J. Park**, E. Amaro, D. Mahajan, B. Thwaites, H. Esmaeilzadeh, "AXGAMES: Towards Crowdsourcing Quality Target Determination in Approximate Computing," in *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, April 2016.
21. H. Sharma, **J. Park**, D. Mahajan, E. Amaro, J. Kim, C. Shao, A. Mishra, H. Esmaeilzadeh "From High-Level Deep Neural Models to FPGAs," in *The 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, October 2016.
22. D. Mahajan, **J. Park**, E. Amaro, H. Sharma, A. Yazdanbaksh, J. Kim, H. Esmaeilzadeh, "TABLA: A Unified Template-based Framework for Accelerating Statistical Machine Learning," in *The 22nd IEEE Symposium on High Performance Computer Architecture (HPCA)*, March 2016.
(Distinguished Paper Award)
23. D. Mahajan, A. Yazdanbaksh, **J. Park**, B. Thwaites, H. Esmaeilzadeh, "Towards Statistical Guarantees in Controlling Quality Tradeoffs in Approximate Acceleration," in *International Symposium on Computer Architecture (ISCA)*, June 2016.
24. A. Yazdanbakhsh, **J. Park**, H. Sharma, P. Lotfi-Kamran, H. Esmaeilzadeh, "Neural Acceleration for GPU Throughput Processors," in *The 48th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, December 2015.
25. **J. Park**, H. Esmaeilzadeh, X. Zhang, M. Naik, W. Harris, "FLEXJAVA: Language Support for Safe and Modular Approximate Programming," in *The 10th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE)*, September 2015.
26. A. Yazdanbakhsh, D. Mahajan, B. Thwaites, **J. Park**, A. Nagendrakumar, S. Sethuraman, K. Ramkrishnan, N. Ravindran, R. Jariwala, A. Rahimi, H. Esmaeilzadeh, K. Bazargan, "AXILOG: Language Support for Approximate Hardware Design," in *Design Automation and Test in Europe (DATE)*, March 2015.
27. R. S. Amant, A. Yazdanbakhsh, **J. Park**, B. Thwaites, H. Esmaeilzadeh, A. Hassibi, L. Ceze, D. Burger, "General-Purpose Code Acceleration with Limited-Precision Analog Computation," in *The 41th International Symposium on Computer Architecture (ISCA)*, June 2014.
(Nominated for CACM Research Highlights; Honorable Mention in IEEE Micro Top Picks)
28. B. Thwaites, G. Pekhimenko, A. Yazdanbakhsh, **J. Park**, G. Mururu, H. Esmaeilzadeh, O. Mutlu, T. Mowry, "Rollback-Free Value Prediction with Approximate Loads," in *The 24th International Conference on Parallel Architectures and Compilation Techniques (PACT)*, August 2014.
29. J. Choi, **J. Park**, J. Seol, and S. Maeng, "Isolated Mini-domain for Trusted Cloud Computing," in *The 13th International Symposium on Cluster, Cloud, and Grid Computing (CCGrid)*, May 2013.
30. **J. Park**, D. Lee, B. Kim, J. Huh, S. Maeng, "Locality-aware Dynamic VM Reconfiguration on MapReduce Clouds," in *The 21st International ACM Symposium on High-Performance Parallel and Distributed Computing (HPDC)*, June 2012.

Refereed Journal Articles

1. H. Ham*, W. Yang*, Y. Shin, O. Woo, G. Heo, S. Lee, **J. Park**, G. Kim, "ONNXim: A Fast, Cycle-level Multi-core NPU Simulator," in *IEEE Computer Architecture Letters (CAL)*, December 2024.
2. S. Moon, S. Hong, M. Kim, D. Seo, J. Kim, R. Ko, S. Choi, J. Cha, J. Kim, S. Lim, H. Lee, H. Park, G. Choi, J. Kim, J. Lee, **J. Park**, J. Kim "LPU: A Latency-optimized and Highly Scalable Processor for Large Language Model Inference" in *IEEE Micro, special issue on Contemporary Industry Products*, 2024.
3. S. Hwang, D. Baek, **J. Park**, J. Huh, "Cerberus: Triple Mode Acceleration of Sparse Matrix and Vector Multiplication," in *IEEE Transactions on Architecture and Code Optimization (TACO)*, 2024.

4. J. Park, S. Kang, S. Lee, T. Kim, **J. Park**, Y. Kwon, and J. Huh, "Hardware Hardened Sandbox Enclaves for Trusted Serverless Computing" in *IEEE Transactions on Architecture and Code Optimization (TACO)*, 2023.
5. S. Noh, J. Koo, S. Lee, **J. Park**, and J. Kung, "FlexBlock: A Flexible DNN Training Accelerator with Multi-Mode Block Floating Point Support" in *IEEE Transactions on Computers (TC)*, 2023.
6. S. Lee, R. Hwang, **J. Park**, and M. Rhu, "HAMMER: Hardware-friendly Approximate Computing for Self-attention with Mean-redistribution and Linearization" in *IEEE Computer Architecture Letters (CAL)*, 2023.
7. W. Seo, S. Cha, Y. Kim, J. Huh, and **J. Park**, "SLO-aware Inference Scheduler for Heterogeneous Processors in Edge Platforms" in *Transactions on Architecture and Code Optimization (TACO)*, 2021.
8. D. Mahajan, K. Ramkrishnan, R. Jariwala, A. Yazdanbakhsh, **J. Park**, B. Thwaites, A. Nagendrakumar, A. Rahimi, H. Esmaeilzadeh, K. Bazargan, "AXILOG: Abstractions for Approximate Hardware Design and Reuse," in *IEEE Micro, special issue on Alternative Computing Designs and Technologies*, October 2015.

Refereed Workshop Papers

1. J. Cho, M. Kim, H. Choi, **J. Park**, "LLMServingSim: A Simulation Infrastructure for LLM Inference Serving Systems", in *ISCA Workshop on ML for Computer Architecture and Systems (MLArchSys)*, June 2024.
2. Y. Lee, **J. Park**, "LVS: A Learned Video Storage for Fast and Efficient Video Understanding" in *Efficient Deep Learning for Computer Vision (ECV) in conjunction with CVPR*, June 2024 (To Appear).
3. H. Sharma, **J. Park**, E. Amaro, B. Thwaites, P. Kotha, A. Gupta, J. Kim, A. Mishra, H. Esmaeilzadeh, "DNNWEAVER: From High-Level Deep Network Models to FPGA Acceleration," in *The Second Workshop on Cognitive Architectures (CogArch) in conjunction with ASPLOS*, April 2016.
4. D. Mahajan, A. Yazdanbakhsh, **J. Park**, B. Thwaites, H. Esmaeilzadeh, "Prediction-Based Quality Control for Approximate Accelerators," in *The Second Workshop on Approximate Computing Across the System Stack (WACAS) in conjunction with ASPLOS*, March 2015.
5. **J. Park**, K. Ni, X. Zhang, H. Esmaeilzadeh, M. Naik, "Expectation-Oriented Framework for Automating Approximate Programming," in *The First Workshop on Approximate Computing Across the System Stack (WACAS) in conjunction with ASPLOS*, March 2014.
6. A. Yazdanbakhsh, B. Thwaites, **J. Park**, H. Esmaeilzadeh, "Methodical Approximate Hardware Design and Reuse," in *The First Workshop on Approximate Computing Across the System Stack (WACAS) in conjunction with ASPLOS*, March 2014.
7. A. Yazdanbakhsh, R. Amant, B. Thwaites, **J. Park**, H. Esmaeilzadeh, A. Hassibi, L. Ceze, D. Burger, "Toward General-Purpose Code Acceleration with Analog Computation," in *The First Workshop on Approximate Computing Across the System Stack (WACAS) in conjunction with ASPLOS*, March 2014.
8. B. Thwaites, A. Yazdanbakhsh, **J. Park**, H. Esmaeilzadeh, "Bio-Accelerators: Bridging Biology and Silicon for General-Purpose Computing," in *Wild and Crazy Ideas (WACI) in conjunction with ASPLOS*, March 2014.

Research Experience

Research Assistant. Alternative Computing Technology (ACT) Lab

Aug. 2013–Aug. 2018

- Georgia Institute of Technology
- Advisor: Prof. Hadi Esmaeilzadeh

Visiting Researcher. Alternative Computing Technology (ACT) Lab

Jan. 2018–Aug. 2018

- University of California, San Diego
- Advisor: Prof. Hadi Esmaeilzadeh

	Research Intern. Architecture Research Group (ARG) <ul style="list-style-type: none"> • NVIDIA Research • Mentors: Dr. Arslan Zulfiqar and Dr. Eiman Ebrahimi • Manager: Dr. Stephen Keckler 	May 2017–Aug. 2017
	Research Intern. Catapult team <ul style="list-style-type: none"> • Microsoft Research • Mentor: Dr. Eric Chung • Manager: Dr. Doug Burger 	Jan. 2016–May 2016
	Research Assistant. Computer Architecture (CA) Lab <ul style="list-style-type: none"> • Korea Advanced Institute of Science and Technology (KAIST) • Advisor: Prof. Seungryoul Maeng 	Feb. 2010–Jul. 2013
Teaching Experience	Instructor. <ul style="list-style-type: none"> • CS230: System Programming • CS610: Parallel Processing • CS311: Computer Organization • CS411: System for Artificial Intelligence • CS510: Computer Architecture • CS230: System Programming • CS311: Computer Organization • CS230: System Programming • CS492: Special Topic in Computer Science: System for Artificial Intelligence • CS230: System Programming • CS492: Special Topic in Computer Science: System for Machine Learning 	Fall 2024 Spring 2024 Spring 2024 Fall 2023 Spring 2023 Fall 2022 Spring 2022 Fall 2021 Spring 2021 Fall 2020 Spring 2020
	Teaching Assistant. <ul style="list-style-type: none"> • CS3220: Processor Design • CS3220: Processor Design • CS8803: Alternative Computing Technology • CS211: Digital System and Lab • CS311: Embedded Computer Systems 	Georgia Institute of Technology Georgia Institute of Technology Georgia Institute of Technology KAIST KAIST
		Fall 2016 Fall 2014 Spring 2014 Spring 2011 Fall 2010
References Available to Contact	Hadi Esmaeilzadeh. Professor, UCSD	hadi@eng.ucsd.edu
	Nam Sung Kim Professor, UIUC	nskim@illinois.edu
	Doug Burger. Technical Fellow and Corporate VP, Microsoft Research	dburger@microsoft.com
	Eric Chung. VP of AI Computing, NVIDIA	eschung@nvidia.com