

<b>CONTACT</b>	poant@nvidia.com	<a href="https://research.nvidia.com/person/po-an-tsai">https://research.nvidia.com/person/po-an-tsai</a>
<b>EDUCATION</b>	<b>Massachusetts Institute of Technology</b> <i>Ph.D.</i> in Computer Science, minor in <i>Optimization Methods</i> <i>S.M.</i> in Computer Science <i>Advisor: Professor Daniel Sanchez</i>	<b>September 2013 – June 2019</b> <b>June 2015</b>
	<b>National Taiwan University</b> <i>B.Sc.</i> in Electrical Engineering	<b>June 2012</b>
<b>SKILLS AND TOOLS</b>	<b>Languages:</b> C, C++, Python, Java, Verilog, Matlab, bash, SQL <b>Libraries:</b> CUDA, OpenCL, matplotlib, Pytorch, TensorFlow <b>Tools:</b> Git, Intel Pin, Timeloop, Zsim, ModelSim, Altera Quantus II, IC Encounter	
<b>WORK EXPERIENCE</b>	<b>NVIDIA Research, Westford MA</b> <i>Sr. Research Scientist</i> <i>Research Scientist</i>	
		<b>July 2022 – Current</b> <b>July 2019 – June 2022</b>
	<p>I develop architectures to address the emerging demands of computer vision and machine learning algorithms. This task requires understanding and analyzing the interplay between hardware, software, and algorithms. Specifically, I work on developing future tensor accelerator that accelerates a wider range of tensor algorithms than conventional accelerators. To evaluate designed accelerators, I use and contribute to an open-source analytical modeling tool (Timeloop+Accelergy) for rapid evaluation of DNN accelerators. My research has influenced how future generations of HW and SW systems will be designed at NVIDIA.</p> <p>I also collaborate with teams across the company, spanning software, research, engineering, and product groups, and publish original research and speak at conferences and events.</p>	
	<b>MIT Computer Science &amp; Artificial Intelligence Lab, Cambridge MA</b> <i>Research Assistant</i>	
		<b>September 2013 – June 2019</b>
	<p>My Ph.D. research focuses on reducing data movement in computer systems to improve their performance and energy efficiency. I designed new memory hierarchies, developed algorithms for data placement and workload scheduling, and co-designed hardware/software to optimize systems.</p> <p>Across my projects, I prototyped ideas extending Zsim, a C++, Intel Pin-based open-sourced multicore simulator. I leveraged latest commodity hardware features (e.g., Intel CAT) and profiled workloads using hardware performance counters. I made essential changes throughout the software stack, including applications (e.g., key-value store, graph analytics) and runtime/compiler in Maxine, a Java-based research JVM.</p>	
	<b>VMware, Palo Alto CA</b> <i>Ph.D. Intern</i>	
		<b>June 2015 – August 2015</b>
	<p>Distributed Resource Management Team. Worked on a VM scheduler that performs multi-dimensional resource balancing and traffic engineering. Proposed a randomized and graph-clustering-based algorithm and evaluated it using a trace-driven simulator written in Python. My algorithm reduces the runtime overhead by 10× while improving utilization by 5% and was publicly released in 2016 and filed as a US patent.</p>	
<b>PATENT</b>	<b>Resource-Based Virtual Computing Instance Scheduling</b> <b>Po-An Tsai, Sahan Gamage, and Rean Griffith.</b>	US 15283274
	<b>Flexible Accelerator for a Tensor Workload</b> <b>Po-An Tsai, Neal Crago, Angshuman Parashar, Joel Springer Emer, Stephen William Keckler</b>	US Patent 17343597,17343582
<b>RECENT PUBLICATIONS</b>	<b>Sparseloop: An Analytical Approach To Sparse Tensor Accelerator Modeling</b> Yannan Nellie Wu, <b>Po-An Tsai</b> , Angshuman Parashar, Vivienne Sze, Joel S Emer.	
		<b>MICRO-55</b>
	<b>SIMD<sup>2</sup>: A Generalized Matrix Instruction Set for Accelerating Tensor Computation beyond GEMM</b> Yunan Zhang, <b>Po-An Tsai</b> , Hung-Wei Tseng.	
		<b>ISCA-49</b>