

EDUCATION	<b>University of Illinois</b> , Urbana Champaign, IL Computer Science Ph.D. Advisor: <a href="#">Prof. Tianyin Xu</a> Start Aug. 2021
	<b>Northwestern University</b> , Evanston, IL M.S. Computer Science, B.S. Electrical Engineering GPA: 4.0/4.0 (Summa Cum Laude) Graduated June 2021
PUBLICATIONS	<ol style="list-style-type: none"> <li>1. <b>[ASPLOS 2022]</b> Brian Suchy, Souradip Ghosh, Drew Kersnar, <b>Siyuan Chai</b>, Zhen Huang, Aaron Nelson, Michael Cuevas, Alex Bernat, Gaurav Chaudhary, Nikos Hardavellas, Simone Campanoni, and Peter Dinda. “<a href="#">CARAT CAKE: replacing paging via compiler/kernel cooperation</a>”. In <i>Proceedings of Architectural Support for Programming Languages and Operating Systems</i>.</li> <li>2. <b>[Radiology]</b> Ramsey M Wehbe, Jiayue Sheng, Shinjan Dutta, <b>Siyuan Chai</b>, Amil Dravid, Semih Barutcu, Yunan Wu, Donald R. Cantrell, Nicholas Xiao, Hatice Savas, Rishi Agrawal, Nishant Parekh, Aggelos K. Katsaggelos. “<a href="#">DeepCOVID-XR: An Artificial Intelligence Algorithm to Detect COVID-19 on Chest Radiographs Trained and Tested on a Large U.S. Clinical Data Set.</a>” <i>Radiological Society of North America</i>.</li> </ol>
RESEARCH EXPERIENCE	<p><b>UIUC Xlab</b>, <a href="#">Prof. Tianyin Xu</a> Aug. 2021 to Present  <i>Support Linux Kernel for Elastic Cuckoo Page Table</i></p> <ul style="list-style-type: none"> <li>• Adapting Linux kernel, primarily the memory management portion, to have a more versatile support for non tree page table designs like <a href="#">Elastic Cuckoo Page Table</a> (ECPT), a hash page table with memory-level parallelism</li> <li>• Analyzed and addressed the adaptation challenges at kernel level including page table management interface, transparent huge pages, and page table isolation</li> <li>• Extensively modified address translation portion of QEMU to simulate ECPT’s hardware behavior</li> </ul> <p><b>NU Compilers Group</b>, <a href="#">Prof. Simone Campanoni</a> Jan. 2021 to July 2021  <i>Enhance Parallelism by Utilizing Commutative Loop Iterations</i></p> <ul style="list-style-type: none"> <li>• Coded a LLVM pass to transform serial code to parallel by telling the commutativity of &lt;load, ALU operation, store&gt; triplet across loop iterations</li> </ul> <p><b>NU Parallelism Group</b>, <a href="#">Prof. Peter Dinda</a> June 2020 to May 2021  <i>CARAT CAKE: Replacing Paging via Compiler/Kernel Cooperation</i></p> <ul style="list-style-type: none"> <li>• Designed and implemented CARAT CAKE, an allocation level address space which aims to replace virtual memory and paging with protection checks inserted at compile time and allocations tracked in runtime</li> <li>• Implemented a competitive paging address space with support for red black tree and splay tree data structures to track VA-PA mapping, transparent huge pages, and PCID; performance measured with performance monitoring counter</li> <li>• Designed runtime protection check with address mapping data structures</li> </ul> <p><b>Image &amp; Video Processing Lab</b>, <a href="#">Prof. Aggelos Katsaggelos</a> June 2019 to July 2021  <i>DeepCOVID-XR</i></p> <ul style="list-style-type: none"> <li>• Designed and implemented a CNN model to flag out positive COVID cases based on patients’ chest X-ray images</li> <li>• Outperformed experienced radiologists with an accuracy of 85% compared to 76 - 82% and AUC of 0.935 compared to 0.819 - 0.856</li> </ul>

WORK EXPERIENCE	<p><b>Software Engineering Intern</b>, Google Cloud Infrastructure May 2022 to Aug. 2022  <i>Machine Model Population Pipeline</i></p> <ul style="list-style-type: none"> <li>Designed a distributed pipeline to collect data of all Google's server machines (4M+) to model their physical topology. It implements batch reads from Bigtable and capacitor or makes RPC calls with rate limitation</li> <li>Validated mac address of machines with as-maintained models across three data sources. Results will be stored in Spanner</li> </ul> <p><b>Research Intern</b>, Tencent Network Group June 2021 to Aug. 2021  <i>Service Driven Network Verification tool</i></p> <ul style="list-style-type: none"> <li>Contributed to design a network verification tool for routing configurations (e.g. BGP, OSPF); it supports quantitative query and covers all data plane with global formal modeling and local simulation</li> <li>Designed easy-to-use geo-based intent language for network verification</li> </ul>
PROJECTS	<p><b>CPU-GPU Simulator for Collaborative Workloads Modeling</b></p> <ul style="list-style-type: none"> <li>Analyzed pros and cons of state-of-art simulators (gem5, gem5-GPU, and UVMSmart) when modeling memory performance of CPU-GPU collaborative workloads</li> <li>Designed and prototyped a CPU-GPU memory subsystem simulator for workloads running on CPU-GPU unified virtual memory.</li> <li>Integrated gem5 and UVMSmart with IPC to model performance of CPU, GPU and on-demand page migration between them</li> </ul> <p><b>Log-Structured Merge-Trees (LSM-trees) Optimization with eBPF</b></p> <ul style="list-style-type: none"> <li>Conducted a comprehensive study of optimizing LSM-trees with eBPF</li> <li>Developed two designs to offload the LSM compaction to eBPF with hook point at FS layer or NVMe driver layer</li> </ul> <p><b>Distributed Machine Learning Inference Framework</b></p> <ul style="list-style-type: none"> <li>Developed a fault-tolerant resource-fair ML inference framework from scratch</li> <li>Coded a totally ordered KV distributed file system with quorum consensus to support the ML model inference and data storage</li> <li>Designed and implemented failure detection with a ring-based membership protocol</li> </ul> <p><b>C-style Language Compiler</b>, CS 322 Compiler Construction</p> <ul style="list-style-type: none"> <li>Created, from scratch, a compiler to translate C-style language to x86_64 assembly</li> <li>Implemented features including graph-coloring register allocation, liveness analysis, instruction selection with tiling, control flow graph, and memory access checking</li> </ul> <p><b>Middle End Analysis for a C-based API</b>, CS 323 Code Analysis &amp; Transformation</p> <ul style="list-style-type: none"> <li>Coded a LLVM pass to reduce calls to a custom C-based API by implementing analysis like reaching-definition, constant propagation and folding, alias analysis, function inlining, and dead code elimination.</li> </ul>
SKILLS	<p><b>Programming languages:</b>  C/C++, Assembly, Python, Java, Go, JavaScript, MATLAB</p> <p><b>System-level Development:</b>  Unix/Linux, QEMU, Docker, GDB, Make, Linker, LLVM, OpenMP</p> <p><b>Artificial Intelligence:</b>  CUDA, PyTorch, Tensorflow, Keras, Image Processing, Computer Vision</p> <p><b>Hardware:</b>  Raspberry Pi, Arduino, VHDL, Verilog</p> <p><b>Web Development:</b>  HTML, CSS, Flask, Django, React</p>
PROFESSIONAL ACTIVITIES	<p><b>OSDI/ATC 2022:</b> Artifact Evaluation Committee  <b>SOSP 2021:</b> Artifact Evaluation Committee, Slack Co-chair</p>