

EDUCATION	<b>University of Illinois</b> , Urbana Champaign, IL Computer Science Ph.D. Candidate Advisor: <a href="#">Prof. Tianyin Xu</a> Research Area: Operating Systems, Memory Systems, SW/HW Codesign 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>[<b>OSDI 2025</b>] <b>Siyuan Chai</b>, Jiyuan Zhang, Jongyul Kim, Alan Wang, Jovan Stojkovic, Weiwei Jia, Dimitrios Skarlatos, Josep Torrellas, and Tianyin Xu. “<a href="#">EMT: An Operating System Framework for New Memory Translation Architectures</a>.” In <i>Proceedings of the 19th USENIX Symposium on Operating Systems Design and Implementation</i>.</li> <li>[<b>ASPLOS 2025</b>] Yan Sun, Jongyul Kim, Douglas Yu, Jiyuan Zhang, <b>Siyuan Chai</b>, Michael Jaemin Kim, Hwayong Nam, Jaehyun Park, Eojin Na, Yifan Yuan, Ren Wang, Jung Ho Ahn, Tianyin Xu, Nam Sung Kim. “<a href="#">M5: Mastering Page Migration and Memory Management for CXL-based Tiered Memory Systems</a>.” In <i>Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems</i>.</li> <li>[<b>ASPLOS 2024</b>] Jiyuan Zhang, Weiwei Jia, <b>Siyuan Chai</b>, Peizhe Liu, Jongyul Kim, and Tianyin Xu. “<a href="#">Direct Memory Translation for Virtualized Cloud</a>” In <i>Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems</i>.</li> <li>[<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 the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems</i>.</li> <li>[<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>
WORK EXPERIENCE	<b>NVIDIA</b> , CUDA Unified Virtual Memory <i>Improving AI Workloads Performance via UVM</i> May 2025 to Present Mentor: Ram Tummala <ul style="list-style-type: none"> <li>Developed user-level caching for UVM, boosting LLM serving performance on UVM by 30.4x overall (81.3x in prefill and 5.25x in decode)</li> </ul> <b>Meta</b> , AI System Co-design, Software Engineering Intern <i>Chunked Prefill</i> May to Aug. 2024 Mentor: <a href="#">Dr. Jaewon Lee</a> <ul style="list-style-type: none"> <li>Prototyped and implemented chunked prefill, a technique mitigates prefill-decode interference in LLM serving by splitting prefill request into smaller chunks</li> <li>Compared to Meta’s production baseline, it offers up to 1.7x better p99 inter-token latency and 1.3x higher serving capacity (max throughput under tail latency constraints).</li> </ul> <b>Google</b> , Google Cloud, Software Engineering Intern <i>Machine Model Population Pipeline</i> May to Aug. 2022 Mentor: <a href="#">Alex Tran</a> <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> </ul>

	<ul style="list-style-type: none"> <li>Validated mac address of machines with as-maintained models across three data sources. Results will be stored in Spanner</li> </ul>	
	<b>Tencent</b> , Network Group, Research Intern <i>Service Driven Network Verification Tool</i>	June to Aug. 2021 Mentor: <a href="#">Dr. Congcong Miao</a>
	<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> </ul>	
RESEARCH EXPERIENCE	<b>UIUC Xlab</b> , <a href="#">Prof. Tianyin Xu</a> <i>EMT: An OS Framework for New Memory Translation Architectures</i>	Aug. 2021 to Present
	<ul style="list-style-type: none"> <li>Designed a hardware-neutral, extensible framework with minimal overhead that supports diverse memory translation schemes(e.g. tree- and hash-based)</li> <li>Built an EMT-Linux-based open platform for prototyping, developing and evaluating memory translation research</li> </ul>	
	<i>Direct Memory Translation for Virtualized Clouds</i>	
	<ul style="list-style-type: none"> <li>Proposed Direct Memory Translation (DMT), a practical hardware-software extension for x86-based address translation; it minimizes address translation overhead by directly fetching PTEs</li> <li>Speeded up page walks by 1.61x and overall application execution by 1.21x in virtualized environment</li> </ul>	
	<b>NU Parallelism Group</b> , <a href="#">Prof. Peter Dinda</a> <i>CARAT CAKE: Replacing Paging via Compiler/Kernel Cooperation</i>	June 2020 to May 2021
	<ul style="list-style-type: none"> <li>Designed and implemented 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 design with support for red black tree and splay tree data structures to track VA-PA mapping, huge pages, and PCID</li> </ul>	
	<b>Image &amp; Video Processing Lab</b> , <a href="#">Prof. Aggelos Katsaggelos</a> <i>DeepCOVID-XR</i>	June 2019 to July 2021
	<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>	
PROJECTS	<b>CPU-GPU Simulator for Collaborative Workloads Modeling</b>	
	<ul style="list-style-type: none"> <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 to model performance of CPU, GPU and on-demand page migration between them</li> </ul>	
	<b>C-style Language Compiler</b>	
	<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>	
	<b>Middle End Analysis for a C-based API</b>	
	<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	<b>Programming languages:</b> C/C++, Assembly, Python, Java, Go, JavaScript, MATLAB	
	<b>Artificial Intelligence:</b> LLM, CUDA, PyTorch, Tensorflow, Keras, Image Processing, Computer Vision	

**System-level Development:**

Linux Kernel, QEMU, Docker, GDB, Nsys, Make, Linker, LLVM, OpenMP

**Hardware:**

GPU, FPGA, Raspberry Pi, Arduino, VHDL, Verilog

**Web Development:**

HTML, CSS, Flask, Django, React

**PROFESSIONAL  
ACTIVITIES**

**ASPLOS 2025:** Artifact Evaluation Committee

**OSDI/ATC 2022, 2023:** Artifact Evaluation Committee

**SOSP 2021:** Artifact Evaluation Committee, Slack Co-chair

**GRANTS**

Travel grants for OSDI 2023 and 2025

**TEACHING  
EXPERIENCE****Teaching Assistant** - University of Illinois Urbana-Champaign

Spring 2025 CS 340: Intro to Computer Systems with [Prof. Luther Tychonievich](#)

Fall 2024 CS 423: Operating System Design with [Prof. Tianyin Xu](#)

Fall 2023 CS 423: Operating System Design with [Prof. Tianyin Xu](#)

Fall 2022 CS 423: Operating System Design with [Prof. Tianyin Xu](#)

Spring 2022 CS 598XU: Reliability of Cloud-Scale Systems with [Prof. Tianyin Xu](#)

**Peer Mentor (Undergraduate TA)** - Northwestern University

Spring 2021 CS 336 - Design & Analysis of Algorithms with [Prof. Jason Hartline](#)

Winter 2021 CS 343 - Operating Systems with [Prof. Peter Dinda](#)

Winter 2020 CS 336 - Design & Analysis of Algorithms with [Prof. Konstantin Makarychev](#)

Fall 2019 CS 336 - Design & Analysis of Algorithms with [Prof. Jason Hartline](#)

Spring 2019 CS 336 - Design & Analysis of Algorithms with [Prof. Jason Hartline](#)

**Teaching Assistant** - Washington University in St. Louis

Spring 2018 ESE 205 Introduction to Engineering Design with [Prof. James Feher](#)