

# Xiang Li

+1 765 637 5476 • [li2068@purdue.edu](mailto:li2068@purdue.edu) • [Google Scholar](#) • [Linkedin Profile](#)

## CAREER SUMMARY

Ph.D. student specializing in the end-to-end optimization of Large Language Model applications. Expertise in designing high-performance RAG pipelines, from novel memory retrieval architectures to low-latency LLM inference. Focused on bridging the gap between complex model architectures and efficient, secure, scalable production deployments.

## EDUCATION

|  |                      |
|--|----------------------|
| <b>Purdue University</b> , West Lafayette  | 2021–2027 (expected) |
| ▪ Ph.D. in Computer Engineering. GPA. 3.87 <ul style="list-style-type: none"><li>• Advisor: Saurabh Bagchi</li><li>• Research Area: System for Machine Learning, Efficient LLM Serving</li></ul> |                      |

|   |           |
|---|-----------|
| <b>Purdue university</b> , West Lafayette | 2015–2019 |
| ▪ B.S. Computer Engineering. GPA. 3.89    |           |

## WORK EXPERIENCE

|  |                          |
|--|--------------------------|
| <b>ML Algorithm Intern</b> - Futurewei Technology Inc. San Jose, CA  | Aug 2025 - Dec 2025      |
| ▪ Developed a novel memory architecture utilizing Graph Neural Networks to address the limited context window and degraded attention result in long-horizon embodied task planning.                        |                          |
| ▪ Designed a dynamic experience graph mechanism to consolidate past agentic experiences with current observations, facilitating more robust state representation and decision-making.                      |                          |
| <b>ML Research Intern</b> - Houston Methodist Research Institute. Houston, TX  | Jun 2023 – Aug 2023      |
| ▪ Built a 3D group equivariant CNN-based segmentation model for non-contrast CT brain scans, reaching 87% soft tissue segmentation accuracy and reducing diagnostic latency for ischemic stroke detection. |                          |
| ▪ Partnered with physicians to translate AI models into clinical pipelines, accelerating imaging-based diagnosis.  |                          |
| <b>Control Software Engineer</b> - Cummins Inc. Columbus, IN   | Jun 2019 – Jul 2021      |
| ▪ Engineered dual-encryption bootloader security for Cummins ECM, eliminating unauthorized firmware tampering across 250k+ vehicles in production.   |                          |
| ▪ Collaborated with cross-functional teams to deliver model-based control software for critical automotive components, ensuring reliability in high-volume manufacturing.                                  |                          |
| ▪ Tuned and validated 30+ OBD parameters for RAM truck compliance with U.S. emission standards, strengthening regulatory approval and system robustness.   |                          |
| <b>Graduate Teaching Assistant</b> -Purdue University. West Lafayette, IN  | Aug 2021 – July 2025     |
| ▪ Teaching assistant for Purdue ECE embedded software senior design course ECE 47700.  |                          |
| ▪ Teaching assistant for Purdue ECE embedded system design course ECE 36200.   |                          |
| ▪ Teaching assistant for Purdue First Year Engineering course ENG 131,132.   |                          |
| <b>Graduate Research Assistant</b> -Purdue University. West Lafayette, IN  | Summer 2022, Summer 2024 |
| ▪ Engineered a unified scheduling system for heterogeneous edge devices and hybrid networks, optimizing resource allocation across diverse computing nodes.  |                          |
| ▪ Established a testbed integrating Private 5G and Wi-Fi 6 to evaluate network performance and reliability under industrial constraints.   |                          |

## RECENT PROJECTS

|  |  |
|--|--|
| <b>Purdue ECE Senior Design Website</b> <a href="#">[link]</a>   |  |
| ▪ Developed an end-to-end web application within the React framework for the instruction team of the senior design course.                     |  |
| ▪ Engineered a modular architecture to synthesize diverse course needs into a unified system, including journaling, inventory, archiving, etc. |  |

## TECHNICAL SKILLS

|   |
|---|
| <b>Languages.</b> C, Python, R, JS ▪ <b>Machine Learning.</b> PyTorch, TensorFlow, LLM Inference Optimization   |
| ▪ <b>H/W Platforms.</b> Jetson AGX, Jetson NX, Jetson Nano, Arduino Uno, Raspberry Pi 4 (baremetal) ▪ <b>Web Technologies.</b> socketio, Wireshark, React |
|   |

|              |  |                                      |
|--------------|--|--------------------------------------|
| COURSEWORK   | <b>Graduate level.</b> Estimation Theory, Random Process, Computational Modeling, Deep Learning, Reinforcement Learning Theory   |                                      |
| PUBLICATIONS | <p>[NIPS@DL4C] <a href="#">Deep-Reproducer: From Paper Understanding to Code Generation</a><br/>           Pengcheng Chen, Ning Yan, Zihan Zhao, Yixiao Lin, Huaibo Chen, Yue Hu, Qinbo Bai, <b>Xiang Li</b>, Masood Mortazavi (2025)</p> <p>[arXiv] <a href="#">Ascendra: Dynamic Request Prioritization for Efficient LLM Serving.</a><br/> <b>Xiang Li*</b>, Azam Ikram*, Sameh Elnikety, Saurabh Bagchi. arXiv preprint (2025)</p> <p>[ESOC] Enhanced Brain Tissue Segmentation In Non-Contrast CT For Ischemic Stroke Diagnosis<br/> <b>Xiang Li</b>, Saurabh Bagchi, John Volpi, Stephen Wong, Kelvin Wong (2025)</p> <p>[arXiv] <a href="#">HopTrack: A Real-time Multi-Object Tracking System for Embedded Devices</a><br/> <b>Xiang Li</b>, Cheng Chen, Yuan-yao Lou, Mustafa Abdallah, Kwang Taik Kim, Saurabh Bagchi. arXiv preprint (2024)</p> <p>[arXiv] <a href="#">Dynamic DAG-application scheduling for multi-tier edge computing in heterogeneous networks</a><br/> <b>Xiang Li</b>, Mustafa Abdallah, Yuan-Yao Lou, Mung Chiang, Kwang Taik Kim, Saurabh Bagchi. arXiv preprint (2023)</p> <p>[SRDS] <a href="#">DAG-based task orchestration for edge computing</a><br/> <b>Xiang Li</b>, Mustafa Abdallah, Shikhar Suryavansh, Mung Chiang, Kwang Taik Kim, Saurabh Bagchi (2022)</p> |                                      |
| ACHIEVEMENTS | <p><b>Tinker Research Grant</b>, Thinking Machine Lab.</p> <ul style="list-style-type: none"> <li>▪ Awarded \$5000 credits for research project uses Tinker.</li> </ul> <p><b>Honor Society Eta Kappa Nu</b>, Purdue University.</p> <ul style="list-style-type: none"> <li>▪ Selected for being in the top 10% of graduate students.</li> </ul> <p><b>Eli Shay Scholarship</b>, Purdue University</p>   | 2025<br>2021<br>2017                 |
| SERVICE      | <p><b>Reviewer</b>, BHI 25'</p> <p><b>Student Administrator</b>, USENIX ATC 24'</p> <p><b>Artifact review committee</b>, OSDI 24'</p> <p><b>Artifact review committee</b>, USENIX ATC 24'</p> <p><b>Reviewer</b>, IEEE Networking Letter</p>   | 2025<br>2024<br>2024<br>2024<br>2023 |