

# Bodun Hu

---

## CONTACT INFORMATION

*E-mail:* bodunhu@utexas.edu  
*Website:* <https://www.bodunhu.com>

2317 Speedway  
The University of Texas at Austin  
Austin, TX 78712 USA

## RESEARCH INTERESTS

Systems for ML, Operating System, heterogeneity, ML SW-HW Co-design, Distributed System

## EDUCATION

### The University of Texas at Austin

Ph.D. in Computer Science  
Advisor: Aditya Akella

### The University of Texas at Austin

M.S. in Computer Science, May 2021  
Advisor: Christopher J. Rossbach

### The University of Texas at Austin

B.S. in Computer Science, May 2020 (Research Distinction)

## PUBLICATIONS

Liyan Tang, Grace Kim, Xinyu Zhao, Thom Lake, Wenxuan Ding, Fangcong Yin, Prasann Singhal, Manya Wadhwa, Zeyu Leo Liu, Zayne Sprague, Ramya Namuduri, **Bodun Hu**, Juan Diego Rodriguez, Puyuan Peng, Greg Durrett. ChartMuseum: Testing Visual Reasoning Capabilities of Large Vision-Language Models. *Preprint*.

**Bodun Hu\***, Luis Pabon\*, Saurabh Agarwawl, Aditya Akella. Patchwork: A Unified Framework for RAG Serving. *Preprint*.

**Bodun Hu**, Shuozhe Li, Saurabh Agarwal, Myungjin Lee, Akshay Jajoo, Jiamin Li, Le Xu, Geon-Woo Kim, Donghyun Kim, Hong Xu, Amy Zhang, Aditya Akella. StitchLLM: Serving LLMs, One Block at a Time. *ACL 25*.

Ajay Jaiswal, **Bodun Hu**, Lu Yin, Yeonju Ro, Shiwei Liu, Tianlong Chen, Aditya Aeklla. FFN-SkipLLM: A Hidden Gem for Autoregressive Decoding with Adaptive Feed Forward Skipping. *EMNLP 24*.

**Bodun Hu**, Le Xu, Jeongyoon Moon, Neeraja J. Yadwadkar, Aditya Akella. MOSEL: Inference Serving Using Dynamic Modality Selection. *EMNLP 24*.

Henrique Fingler, Isha Tarte, Hangchen Yu, Ariel Szekely, **Bodun Hu**, Aditya Akella, Christopher J. Rossbach. Towards a Machine Learning-Assisted Kernel with LAKE. *ASPLOS 23*.

**Bodun Hu** and Christopher J. Rossbach. 2020. Altis: Modernizing GPGPU Benchmarks. *ISPASS 20*.

## RESEARCH EXPERIENCE

**Meta**, Menlo Park, CA, USA.

*Student Researcher*

**2025-Current**

Designed and optimized high-performance network stack to reduce latency in large language model (LLM) training.

**Meta**, Menlo Park, CA, USA.

	<p><i>PhD SWE Intern</i> <span style="float: right;"><b>2025</b></span>  Implement efficient communication collectives for LLM and recommendation models.</p> <p><b>The University of Texas at Austin</b> (UT Austin), Austin, TX, USA.</p> <p><i>Research Assistant</i> <span style="float: right;"><b>2021 - Current</b></span>  Implement efficient inference serving systems for LLMs and GenAI.</p> <p><b>Intel</b>, San Jose, CA, USA.</p> <p><i>P4 Dataplane Intern</i> <span style="float: right;"><b>2022</b></span>  <i>TCP-INT</i>: Improved Network Telemetry in TCP Transport for better e2e visibility and improved closed-loop control of TCP workloads.</p> <p><b>The University of Texas at Austin</b> (UT Austin), Austin, TX, USA.</p> <p><i>Research Assistant</i> <span style="float: right;"><b>2017 - 2021</b></span>  <i>LAKE</i>: Built a generic API remoting system to expose accelerator APIs to OS kernel with close-to-native performances.  <i>ALTIS</i>: Designed a benchmark with improved diversity over existing GPU benchmarks by extending application domains with modern CUDA features.</p> <p><b>The University of Texas at Austin</b> (UT Austin), Austin, TX, USA.</p> <p><i>Research Assistant</i> <span style="float: right;"><b>2020</b></span>  <i>TAS</i>: Ported TAS into P4 to facilitate TCP fast-path migration to programmable NICs.</p> <p><b>The University of Texas at Austin</b> (UT Austin), Austin, TX, USA.</p> <p><i>Research Assistant</i> <span style="float: right;"><b>2016 - 2017</b></span>  <i>G-Code-gen</i>: Designed an automated detection system utilizing readily available hardware, which detects and terminates 3D printing processes upon identification of object defects.</p>
INDUSTRY EXPERIENCE	<p><b>H3C</b>, Chengdu, China.</p> <p><i>Software Engineering Intern</i> <span style="float: right;"><b>2018</b></span>  Devised and implemented a highly effective caching strategy, resulting in a significant reduction of video streaming processing latency on Kubernetes cluster by a factor of 3x.</p> <p><b>Wisesoft</b>, Chengdu, China.</p> <p><i>Software Engineering Intern</i> <span style="float: right;"><b>2017</b></span>  Developed a data preprocessing pipeline for improved audio classification in an air traffic control system.</p>
HONORS AND AWARDS	<p>Wesley W. Calhoun Jr. Endowed Scholarship</p> <p>ISPASS Student Travel Award, 2020</p> <p>Research Distinction by the College of Natural Sciences (UT Austin), 2020.</p>
TEACHING	<p><b>CS395T: Advanced Topics in Systems and GenAI (graduate)</b>  Teaching Assistant, UT Austin, Fall 2025</p> <p><b>CS378: System For Machine Learning and Big Data (undergrad)</b>  Teaching Assistant, UT Austin, Fall 2024</p> <p><b>CS378: Multicore Operating System Implementation (undergrad)</b>  Teaching Assistant, UT Austin, Spring 2020</p>

## TALKS

- *Altis: Modernizing GPGPU Benchmarking*, ISPASS'20 (August 2020)
- *Accelerating Kernel Access to Hardware Acceleration*, Texas Systems Symposium (November 2020)

## SERVICE

- Reviewer for ACL 2025, ACL-SRW 2025
- Junior Graduate Admissions Committee, UT Austin (January 2021)