

# Shengze (William) Wang

Ph.D. Candidate, Computer Science and Engineering

Homepage: <https://shengze.io>

Email: shengze@ucsc.edu

## TECHNICAL SKILLS

- **Academic:** Computer Networks & Distributed Systems and their applications in AI Infrastructure (Systems & Networking for AI), AI for Systems & Databases, LLM Inference & Serving Efficiency, VectorDB
- **Languages & Databases:** C/C++, Rust, Python, eBPF, Redis, Cassandra, DynamoDB, Milvus, Pinecone, HDFS, SQL
- **Platforms & Tools:** Linux, ROS, Ray, Kubernetes - k8s, Apache Spark, Kafka, ns-3, P4 Switch, RDMA, CUDA, FPGAs, vLLM, SGLang, JAX/XLA, PyTorch, Transformers, NCCL, TPU, FPGAs, AWS, Azure, GCP, Docker, Slurm, Git
- **Highlights:** Scalable Key-Value Storage, LLM Inference & ML System Design, Applied AI for Systems & Databases, MLOps & Production Deployment, Advanced Network Protocols & Security, Workload Characterization & Balancing, Content Delivery & High-Performance Networks, ML Stream Processing & Real-Time Analytics, Hashing Algorithms, Edge-Cloud Systems

## EDUCATION

<b>University of California, Santa Cruz (UCSC)</b>	San Francisco Bay Area
• <i>Ph.D. in Computer Science and Engineering; Regents Fellowship; BE Dean's Fellowship</i>	2023 - present
<i>Relevant Courses: Adv Computer Networks, Network Security, Computer Architecture, Adv Distributed Systems, Stream Processing, Programming Languages, Adv Machine Learning, Adv Natural Language Processing, Artificial Intelligence, Analysis of Algorithms</i>	
<b>University of North Texas (UNT)</b>	Dallas - Fort Worth
• <i>B.S. in Computer Science; GPA: 4.0; Outstanding Award (Top 1 of Class 2023); President's List</i>	2020 - 2023
<i>Relevant Courses: Algorithms, Machine Learning, Software Engineering, Systems Programming, Database Systems, Computer Networks, Computer Security, Operating Systems, Probability Models, Linear Algebra, IT Project Management, Technical Writing</i>	
<b>King's College London (KCL)</b>	London, United Kingdom
• <i>Visiting Student in Computer Science; Scored: 95/100; JEISE Scholarship</i>	2019

## WORKING EXPERIENCE

• <b>Software Engineering Intern - AI Networking R&amp;D</b>	Google LLC	2026 -
• <b>Graduate Student Researcher</b>	Baskin School of Engineering, UCSC	June. 2023 - present
• <b>Undergraduate Research Assistant</b>	Department of Computer Science, UNT	Sept. 2021 - May. 2023

## SELECTED PROJECTS

<b>Scalable AI Infrastructure for High-Performance LLM Serving — C/C++, Python</b>	2023 - present
• <i>Qian Lab, <a href="https://users.soe.ucsc.edu/~qian/">https://users.soe.ucsc.edu/~qian/</a>; Center for Research in Systems and Storage (CRSS)</i>	
<ul style="list-style-type: none"><li>○ Designed methods to <b>provide user anonymity and low-overhead encryption for queries/responses in large-scale overlay networks</b>, blending failover resilience with robust privacy guarantees. Developed <b>load balancing mechanisms with distributed key management</b> to support a decentralized and fault-tolerant LLM system.</li><li>○ Integrate Vector Databases to <b>accelerate retrieval-augmented generation over billions of embeddings</b>, achieving efficient approximate nearest neighbor lookups for knowledge-intensive LLM applications. Investigate novel indexing and query pipelines to enhance retrieval accuracy and throughput.</li><li>○ Devised a <b>draft-then-filter</b> mechanism to generate candidate tokens, then selectively offload low-confidence tokens to the full-scale LLM, achieving up to <b>2x speedup in real-time inference</b> without sacrificing output quality.</li><li>○ Implemented <b>NLL-based confidence scoring</b> to dynamically filter high-quality drafts, preventing unnecessary requests to the target LLM and thus lowering GPU consumption. Enhanced MLOps with continuous integration and robust monitoring, <b>streamlining the model lifecycle from draft-model updates to large-scale deployments</b>.</li><li>○ Proposed <b>“CALID”</b>: a novel inference framework that integrates <b>speculative decoding</b> principles to <b>boost throughput and reduce computational overhead</b> for large language models.</li><li>○ Proposed <b>Span-Level Fine-Tuning with unlikelihood training</b>: a novel approach that leverages annotated unfaithful spans in LLM-generated summaries to <b>reduce hallucinations and improve factual accuracy</b>.</li><li>○ Proposed <b>“GenTorrent”</b>: a decentralized overlay network to <b>enhance Large Language Model (LLM) serving scalability and cost-efficiency</b> by harnessing computing resources from distributed contributors. It addresses fundamental challenges in decentralized LLM serving, including <b>overlay network organization, anonymous communication for privacy, efficient overlay forwarding for load balancing and cache reuse, and decentralized verification of model serving quality</b>. GenTorrent aims to <b>democratize AI innovation, significantly reduce serving latency, and improve user privacy</b>, offering a novel approach to future AI deployment.</li></ul>	

<b>Resource Storage and Discovery in Network &amp; Database Systems — C/C++, Rust</b>	2023 - present
• <i>Qian Lab, <a href="https://users.soe.ucsc.edu/~qian/">https://users.soe.ucsc.edu/~qian/</a>; NSF Center for Systems and Storage, <a href="https://ssrc.us/">https://ssrc.us/</a></i>	
<ul style="list-style-type: none"><li>○ Investigate fundamental problems in emerging networks, emphasizing <b>efficient data placement, fault tolerance, and high-throughput designs</b>. (e.g., datacenter networks, CDNs, and quantum networking)</li><li>○ Architect and refine critical components—network protocols, routing algorithms, hashing strategies, and load balancers—for <b>enterprise-scale deployments</b> (e.g., HPC clusters, IoT networks, programmable switches).</li><li>○ Implement and evaluate prototypes using <b>event-driven simulators</b> (e.g., ns-3, p2psim) and cloud-based testbeds (AWS, Lambda Labs, CloudLabs, Supercomputers), leveraging asynchronous I/O and concurrency.</li></ul>	

- Proposed **“LEAD”**: A novel Distributed Learned Hash Table that embeds machine-learned models within Distributed hash table structures to **significantly optimize range query performance** for distributed networked systems. LEAD outperforms existing range-query solutions by demonstrating **superior scalability, reduced latency, and robustness against network churns**. LEAD opens a completely new field for further research on integrating learned models with distributed systems.(<https://github.com/ShengzeWang/LEAD>; <https://github.com/ShengzeWang/RM>)
- Proposed **“Vortex”**: A fully decentralized, planet-scale Vector overlay. Designed Distributed Learned Hashing (DLH) for locality and load balance, DHT routing for fault-tolerant lookup under churn, and per-peer D-HNSW for high-recall local search. Results match SOTA centralized systems' accuracy/latency while reducing per-peer index memory by  $\sim 100\times$ . (Learned Hash Function Library: <https://github.com/ShengzeWang/LearnedHash>)

### Vehicular Edge Computing and Connected Autonomous Vehicles — Python, ROS

2021 - 2023

- *NSF Center for Electric, Connected and Autonomous Technologies*, <https://ecat.center/>, <http://veclab.org/>

- **Profiled hardware resource usage (GPU, CPU, Memory)** for real-time object detectors (YOLO, Faster R-CNN, SSD) deployed in ROS-based CAV perception pipelines.
- Investigated model optimizations, including **quantization, architectural pruning, and mixed-precision**, achieving a measurable trade-off between inference speed and detection accuracy under edge-device constraints.
- Characterized **memory contention and identified performance bottlenecks**, enabling targeted optimization strategies (e.g., improved scheduling, memory partitioning) that enhanced detection throughput.
- Developed **workload models** reflecting diverse edge-device configurations (e.g., NVIDIA Jetson, Intel CPUs, Raspberry Pi), facilitating **informed resource allocation and adaptive scheduling** across heterogeneous deployments.
- Implemented and validated **Vehicle-to-Edge (V2X) communication frameworks** using AWS Edge Services, resulting in reduced latency and improved real-time responsiveness

### False Discovery Rates (FDR) Control in Metaproteomics Search — C++, Python

2021 - 2023

- *Center for Computational Epidemiology and Response Analysis (CeCERA)*, <https://cerl.unt.edu/>

- Addressed systematic FDR biases in metaproteomics pipelines by incorporating **probabilistic modeling and statistical corrections**, reducing false-positive identifications across large proteomic datasets.
- Proposed **“FineFDR”**: an open-source, fine-grained FDR assessment framework that seamlessly integrates with Comet and Percolator outputs at multiple taxonomic ranks. (<https://github.com/Biocomputing-Research-Group/FDR>)
- Implemented the **Expectation-Maximization General-Mixture Model** for clustering proteomic samples based on abundance profiles, substantially enhancing the detection sensitivity for lower-abundance peptides.
- Benchmarked six FDR control solutions (including Comet, Percolator, and Tailor) on ten diverse datasets, demonstrating notable gains in **precision and increased peptide/protein identifications** compared to state-of-the-art approaches.

## SELECTED PUBLICATIONS

• <b>A Distributed Learned Hash Table</b>	Feb. 2026
• <i>IEEE/ACM Transactions on Networking (TON)</i>	<i>First Author</i>
• <b>PlanetServe: A Decentralized, Scalable, and Privacy-Preserving Overlay for Democratizing Large Language Model Serving</b>	Sep. 2025
• <i>2026 USENIX Symposium on Networked Systems Design and Implementation (NSDI)</i>	<i>Co-primary Author</i>
• <b>LEAD: A Distributed Learned Hash Table</b>	Aug. 2025
• <i>2025 IEEE International Conference on Network Protocols (ICNP)</i>	<i>First Author &amp; Oral</i>
• <b>Vortex: Efficient Decentralized Vector Overlay for Similarity Search and Delivery</b>	Aug. 2025
• <i>2025 IEEE International Conference on Network Protocols (ICNP); *BEST POSTER AWARD*</i>	<i>First Author</i>
• <b>Characterizing Perception Deep Learning Algorithms and Applications for Vehicular Edge Computing</b>	Jan. 2025
• <i>Algorithms 2025, 18(1), 31; Special Issue: Machine Learning for Edge Computing</i>	<i>Co-Author</i>
• <b>CALID: Collaborative Accelerate LLM Inference with Draft Model with Filter Decoding</b>	May. 2024
• <i>Poster at 2024 BayLearn - Machine Learning Symposium (Apple)</i>	<i>Co-Author</i>
• <b>Enhancing Faithfulness in Abstractive Summarization via Span-Level Fine-Tuning</b>	May. 2024
• <i>Poster at 2024 BayLearn - Machine Learning Symposium (Apple)</i>	<i>Co-Author</i>
• <b>Distributed Learned Hash Table</b>	Sept. 2024
• <i>2024 IEEE International Conference on Network Protocols (ICNP)</i>	<i>First Author &amp; Poster</i>
• <b>Perception Workload Characterization and Prediction on the Vehicular Edges</b>	Jul. 2023
• <i>2023 IEEE International Conference on Edge Computing (EDGE)</i>	<i>Co-primary Author</i>
• <b>Fine-grained Taxonomy-specific False Discovery Rates Control in Metaproteomics</b>	Nov. 2022
• <i>2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)</i>	<i>First Author &amp; Oral</i>

## PROFESSIONAL SERVICES

• <b>Reviewer</b>	IEEE/ACM TON, IEEE TDSC, ACM SIGCOMM, IEEE INFOCOM
• <b>Teaching Assistant</b>	CSE 13S: Computer Systems and C Programming - 24 Winter, 25 Winter
• <b>Mentor</b>	NSF Research Experiences for Undergraduates (REU) in Vehicular Edge Computing and Security