# YANCHEN LIU

University of Southern California, Los Angeles, CA, USA

+1 (213) 921-6117 | liuyanch@usc.edu | https://jamesnulliu.github.io

## EDUCATION HISTORY

**M.Sc. in Computer Science**                                      **Sep. 2025 – Jun. 2027**
University of California, LA, CA, USA

**B.Eng. in Computer Science**                                     **Sep. 2021 – Jun. 2025**
Shanghai University, Shanghai, China

## HONORS AND AWARDS

1. [2024] _First Prize_ and _Group Competition Award_ in 2024 ASC Student Supercomputer Challenge Global Final.
2. [2022] _First-Class Academic Scholarship_ for outstanding academic performance, Shanghai University.

## PROFESSIONAL EXPERIENCE

**Graduate Research Intern**                                       **Jun. 2025 - Present**
_University of Southern California | INK Lab_
**Topic**: RLHF and Reasoning for LLMs

1. Enhanced _Chain-of-Thought_ data with _segment and token importance_ evaluated by different methods; Utilized the data for _Supervised Fine-Tuning_ of various reasoning models with _TRL_.
2. Modified _v1 engine of vllm_ to implement and reproduce _soft-thinking_ and _latent-thinking_ for reasoning models.
3. Modified _verl_ to selectively _filter out target tokens_ (generated during inference process) based on _entropy aggregation_, which are later used for model updating.

**Machine Learning Engineer Intern**                              **Jul. 2024 - Jun. 2025**
_Shanghai AI Laboratory_
**Topic**: LLM Inference Engine, AI Compiler and Model Fine-tuning

1. Extended _vllm_ for _LLM inference_ on in-house TPUs, including researching and adapting _Speculative Decoding_, _Paged-Attention_ and _Continuous-Batching_.
2. Developed _high-performance kernels_ with _MLIR_ to ensure seamless compatibility and optimal performance of LLMs on in-house TPUs.
3. Conducted research on _LLM knowledge injection and fine-tuning_ for _kernel fusion and translation_ across different hardware platforms.

**Undergraduate Research Assistant**                              **Mar. 2023 - Apr. 2025**
_Shanghai University | Shanghai University Cyber Security Lab_
**Topic**: Vehicle Modeling, Simulation, and Intrusion Detection

1. Combined _deep learning_ and _traditional math modeling_ to assess cyber security and functional safety in _in-vehicle communication system_.
2. Developed a novel efficient _gradient descent solver_ for _Multi-Dimensional Hawkes Process_, and implemented algebraic simplification, vectorization and JIT compilation for optimization; Introduced a novel _MDHP-LSTM structure_ for improved feature extraction in in-vehicle communication data of ECUs and related applications.

**Team Leader of Shanghai University Super-Computing Team**        **Sep. 2023 - Jul. 2024**
_Shanghai University | SHUSCT_
**Topic**: Super-Computing

1. Participated in _2024 ASC Student Supercomputer Challenge_.
2. Assembled, benchmarked and optimized a _high-performance computing cluster_ for running LLMs and super-computing programs.
3. Developed _a custom LLM inference engine_ with various parallelism and scheduling policies; Integrated FlashAttention, quantization and pruning techniques to speed up inference.
4. Profiled and improved performance of _various super-computing programs_ by applying vector instructions, OpenMP, loop unrolling, and MPI.

**Undergraduate Research Intern**                                              **Jun.2023 – Jul. 2024**

*Shanghai University | East China Air Traffic Control Bureau*

**Topic**: Super Resolution for Meteorological Data

1. Deployed cutting-edge deep learning models including *PanGu* and *FourCastNet* to improve *meteorological monitoring, forecasting, and super-resolution* for the East China Air Traffic Control Bureau, replacing traditional numerical methods.

2. Utilized *key-point and semantic constraints*, combined with *feature map fusion and redesigned loss functions*, to improve the up-sampling process in *super-resolution* models, addressing the issue of semantic detail loss in certain areas.

## PUBLICATIONS

1. [2025 | Preprint] J. Lv, X. He, **Y. Liu**, A. Shen, X. Dai, Y. Li, J. Hao, J. Ding, Y. Hu, S. Yin. "*HPCTransEval: A Benchmark of High-Performance GPU-to-CPU Transpilation with Pre-trained Large Language Models*".

2. [2025 | Preprint | Code] Q. Liu†, **Y. Liu†**, R. Li, C. Cao, Y. Li*, X. Li*, P. Wang, R. Feng, "*MDHP-Net: Detecting an Emerging Time-Exciting Threat in IVN*".

3. [2025 | Preprint] Z. Xu, A. Shen, D. Kong, X. Dai, J. Liu, **Y. Liu**, L. Wang, S. Wei, Y. Hu and S. Yin*, "*LLMEngine: Disaggregated Mapping and Memory Management Co-scheduling for Wafer-scale Chips*".

4. [2024 | IEEE Internet of Things Journal | Code] Q. Liu, X. Li, K. Sun, Y. Li* and **Y. Liu***, "*SISSA: Real-Time Monitoring of Hardware Functional Safety and Cybersecurity With In-Vehicle SOME/IP Ethernet Traffic*".

5. [2024 | MDPI Future Internet | Code] Li, X., R. Li, and **Y. Liu**. "*HP-LSTM: Hawkes Process–LSTM-Based Detection of DDoS Attack for In-Vehicle Network*".