

# Xinyu Yang

Carnegie Mellon University

✉ xinyuya2@andrew.cmu.edu | 🏠 xinyuyang.me/ | 📱 Hanyuezhuhua | 🎓 Xinyu Yang

## Education

### Carnegie Mellon University

Ph.D. student in Electrical and Computer Engineering

Pittsburgh, US

Aug. 2023 - Present

### Shanghai Jiao Tong University

Bachelor of Engineering in Computer Science, ACM Honors Class

Shanghai, China

Sept. 2019 - Jun. 2023

- **ACM Honors Class** is an elite CS program for students ranked in the top 5% with aspirations in research.
- Overall GPA: **91.26 / 100**, Ranking: **1 / 29**

## Research Interests

Machine Learning and System, Scalable and Generalizable Foundation Models for In-the-wild Applications

## Publications

\* indicates equal contributions

### APE: Faster and Longer Context-Augmented Generation via Adaptive Parallel Encoding

Xinyu Yang, Tianqi Chen, Beidi Chen

Under Review

### VcLLM: Video Codecs are Secretly Tensor Codecs

Ceyu Xu\*, Yongji Wu\*, Xinyu Yang\*, Beidi Chen, Matthew Lentz, Danyang Zhuo, Lisa Wu Wills

Preprint

### S<sup>2</sup>FT: Efficient, Scalable and Generalizable LLM Fine-tuning by Structured Sparsity

Xinyu Yang, Jixuan Leng, Geyang Guo, Jiawei Zhao, Ryumei Nakada, Linjun Zhang, Huaxiu Yao, Beidi Chen

NeurIPS 2024

### It Takes Two: On the Seamlessness between Reward and Policy Model in RLHF

Taiming Lu\*, Lingfeng Shen\*, Xinyu Yang\*, Weiting Tan, Beidi Chen, Huaxiu Yao

Under Review

### Blindness and Hallucinations: Revisiting Multi-modal Alignment in Vision-Language Large Models

Xinyu Yang, Chenhang Cui, Jaehong Yoon, Yiyang Zhou, Yi-Lin Sung, Mohit Bansal, Beidi Chen, Huaxiu Yao

Under Review

### Zeroth-Order Fine-Tuning of LLMs with Extreme Sparsity

Wentao Guo, Jikai Long, Yimeng Zeng, Zirui Liu, Xinyu Yang, Yide Ran, et al.

Preprint

### TriForce: Lossless Acceleration of Long Sequence Generation with Hierarchical Speculative Decoding

Hanshi Sun, Zhuoming Chen, Xinyu Yang, Yuandong Tian, Beidi Chen

COLM 2024

### Get More with LESS: Synthesizing Recurrence with KV Cache Compression for Efficient LLM Inference

Harry Dong, Xinyu Yang, Zhenyu Zhang, Zhangyang Wang, Yuejie Chi, Beidi Chen

ICML 2024

### Holistic analysis of hallucination in gpt-4v (ision): Bias and interference challenges

Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu, Linjun Zhang, James Zou, Huaxiu Yao

Preprint

### Improving Out-of-Domain Generalization with Domain Relations

Huaxiu Yao\*, Xinyu Yang\*, Xinyi Pan, Shengchao Liu, Pang Wei Koh and Chelsea Finn

ICLR 2024 (Spotlight)

### Multi-domain Long-Tailed Learning By Augmenting Disentangled Representations

Xinyu Yang\*, Huaxiu Yao\*, Allan Zhou and Chelsea Finn

TMLR

### FlatFormer: Flattened Window Attention for Efficient Point Cloud Transformer

Zhijian Liu\*, Xinyu Yang\*, Haotian Tang, Shang Yang and Song Han

CVPR 2023

### BEVFusion: Multi-Task Multi-Sensor Fusion with Unified Bird's-Eye View Representation

Zhijian Liu\*, Haotian Tang\*, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela Rus and Song Han

ICRA 2023

### Variational Inference for Training Graph Neural Networks in Low-Data Regime through Joint Structure-Label Estimation

Danning Lao\*, Xinyu Yang\*, Qitian Wu and Junchi Yan

KDD 2022