

CONTACT INFORMATION	School of Computer Science, Shanghai Jiao Tong University, 800 Dongchuan Road, Shanghai, China Homepage: <a href="https://zhitengli.github.io/">https://zhitengli.github.io/</a>	ieeezhitengli@gmail.com Tel: +86-134-140-18498 Google Scholar, Github
RESEARCH INTERESTS	LLM pretraining and scaling; Image/Video generation; Model compression and inference acceleration (e.g., Quantization); Data synthesis.	
EDUCATION	<b>Shanghai Jiao Tong University</b> , Shanghai, China M.S. Student, Computer Science and Technology • GPA: 3.94/4.0 (Top 5%)	Sep 2023 – Present Advisor: <i>Prof. Yulun Zhang</i>
	<b>Shanghai Jiao Tong University</b> , Shanghai, China B.S., Computer Science and Technology (IEEE Honor Class) • GPA: 4.0/4.3 (Rank 1/113)	Sep 2019 – Jun 2023
PROFESSIONAL EXPERIENCE	<b>ByteDance Seed, Beijing, China</b> Research Intern Project: LLM pretraining with MuP & Scaling Law.	May 2025 – Present Mentors: <i>Huizhuo Yuan, Prof. Quanquan Gu</i>
	<b>Xiaohongshu, Shanghai, China</b> Research Intern Project: Video DiT acceleration (applied our QuantCache to real-world business scenarios).	Jan 2025 – Apr 2025 Mentor: <i>Shuang Sun</i>
	<b>Sony AI, Remote</b> Research Intern Project: Generative Data Augmentation (GenDataAgent is accepted at <b>ICLR 2025</b> ).	Dec 2023 – May 2024 Mentor: <i>Lele Chen</i>
	<b>Amazon Shanghai AI Lab (ASAILab), Shanghai, China</b> Applied Scientist Intern Project: Contributed to <b>DGL</b> & Graph Transformer research.	Jun 2022 – Sep 2023 Mentors: <i>Mufei Li, Minjie Wang</i>
PUBLICATIONS	5×ICLR, 1×ICML, 1×ICCV, 8×Preprints. * indicates equal contribution.	
	<ol style="list-style-type: none"> <li>1. <b>[ICLR 2025] ARB-LLM: Alternating Refined Binarizations for Large Language Models.</b> <i>Zhiteng Li*</i>, Xianglong Yan*, Tianao Zhang, Haotong Qin, Dong Xie, Jiang Tian, Zhongchao Shi, Linghe Kong, Yulun Zhang, and Xiaokang Yang.</li> <li>2. <b>[ICLR 2025] GenDataAgent: On-the-fly Dataset Augmentation with Synthetic Data.</b> <i>Zhiteng Li</i>, Lele Chen, Jerone Andrews, Yunhao Ba, Yulun Zhang, and Alice Xiang.</li> <li>3. <b>[ICLR 2026] DVD-Quant: Data-free Video Diffusion Transformers Quantization.</b> <i>Zhiteng Li*</i>, Hanxuan Li*, Junyi Wu, Kai Liu, Haotong Qin, Linghe Kong, Guihai Chen, Yulun Zhang, and Xiaokang Yang.</li> <li>4. <b>[ICCV 2025] QuantCache: Adaptive Importance-Guided Quantization with Hierarchical Latent and Layer Caching for Video Generation.</b> <i>Junyi Wu*</i>, <i>Zhiteng Li*</i>, Zheng Hui, Linghe Kong, Yulun Zhang, and Xiaokang Yang.</li> <li>5. <b>[ICLR 2026] Quant-dLLM: Post-Training Extreme Low-Bit Quantization for Diffusion Large Language Models.</b> <i>Tianao Zhang*</i>, <i>Zhiteng Li*</i>, Xianglong Yan, Haotong Qin, Yong Guo, and Yulun Zhang.</li> <li>6. <b>[ICLR 2026] PT<sup>2</sup>-LLM: Post-Training Ternarization for Large Language Models.</b> <i>Xianglong Yan*</i>, Chengzhu Bao*, <i>Zhiteng Li</i>, Tianao Zhang, Kaicheng Yang, Haotong Qin, Ruobing Xie, Xingwu Sun, and Yulun Zhang.</li> </ol>	

7. [ICML 2025] **BiMaCoSR: Binary One-Step Diffusion Model Leveraging Flexible Matrix Compression for Real Super-Resolution.**  
 Kai Liu, Kaicheng Yang, Zheng Chen, **Zhiteng Li**, Yong Guo, Wenbo Li, Linghe Kong, and Yulun Zhang.

PREPRINTS

1. **AdaSVD: Adaptive Singular Value Decomposition for Large Language Models.**  
**Zhiteng Li\***, Mingyuan Xia\*, Jingyuan Zhang, Zheng Hui, Haotong Qin, Linghe Kong, Yulun Zhang, and Xiaokang Yang.
2. **BinaryHPE: 3D Human Pose and Shape Estimation via Binarization.**  
**Zhiteng Li**, Yulun Zhang, Jing Lin, Haotong Qin, Jinjin Gu, Xin Yuan, Linghe Kong, and Xiaokang Yang.
3. **VEQ: Modality-Adaptive Quantization for MoE Vision-Language Models.**  
 Guangshuo Qin\*, **Zhiteng Li\***, Zheng Chen, Weihang Zhang, Linghe Kong, and Yulun Zhang.
4. **ReCalKV: Low-Rank KV Cache Compression via Head Reordering and Offline Calibration.**  
 Xianglong Yan\*, **Zhiteng Li\***, Tianao Zhang, Haotong Qin, Linghe Kong, Yulun Zhang, and Xiaokang Yang.
5. **D<sup>2</sup>Quant: Accurate Low-bit Post-Training Weight Quantization for LLMs.**  
 Xianglong Yan\*, Chengzhu Bao\*, **Zhiteng Li**, Tianao Zhang, Haotong Qin, Shaoqiu Zhang, Ruobing Xie, Xingwu Sun, and Yulun Zhang.
6. **Progressive Binarization with Semi-Structured Pruning for LLMs.**  
 Xianglong Yan, Tianao Zhang, **Zhiteng Li**, Haotong Qin, and Yulun Zhang.
7. **CondiQuant: Condition Number Based Low-Bit Quantization for Image Super-Resolution.**  
 Kai Liu, Dehui Wang, **Zhiteng Li**, Zheng Chen, Yong Guo, Wenbo Li, Linghe Kong, and Yulun Zhang.
8. **Low-bit model quantization for deep neural networks: A survey.**  
 Kai Liu, Qian Zheng, Kaiwen Tao, **Zhiteng Li** et al.

ACADEMIC  
SERVICE

- Reviewer
- International Conference on Machine Learning (ICML) 2025, 2026
  - Advances in Neural Information Processing Systems (NeurIPS) 2025
  - International Conference on Learning Representations (ICLR) 2025, 2026
  - International Conference on Computer Vision (ICCV) 2025
  - Computer Vision and Pattern Recognition (CVPR) 2025, 2026

HONORS AND  
AWARDS

- Yang Yuanqing Education Foundation 2025
- Excellent Graduate of Shanghai Jiao Tong University 2023
- Ruiyuan-Sequoia Talent Development Fund 2022
- National Scholarship for Undergraduate Excellence (Top 0.2% nationwide) 2021
- Excellent Undergraduate Scholarship 2020, 2021, 2022

SKILLS

- Programming: Python, Pytorch, C/C++, L<sup>A</sup>T<sub>E</sub>X