

# Lingyun Yang (杨凌云)

Updated June 22, 2025

**Email:** [lyangbk@cse.ust.hk](mailto:lyangbk@cse.ust.hk)

**GitHub:** [mental2008](https://github.com/mental2008)

**LinkedIn:** [stephenyang1999](https://www.linkedin.com/in/stephenyang1999)

**Phone:** (+86) 135-0284-6103

**Office:** BDI 101, UC, HKUST

**Web:** <https://www.lingyunyang.com/>

**RESEARCH INTERESTS** I have a broad interest in resource management for large-scale data centers. Specifically, my research focuses on: (a) improving *resource efficiency* for AI/GPU clusters; (b) building *efficient* and *low-cost* AI model serving systems. My first-authored papers have been published in top-tier systems conferences (NSDI, ATC, SoCC).

**PROFESSIONAL EXPERIENCE** **Alibaba Group** Hangzhou, China  
*Alibaba Holding–Aicheng Technology–Tech Infra and Reliability Engineering (TRE)*  
Senior Engineer (P7) TBA  
◊ Director: Dr. Yinghao Yu

**EDUCATION** **Hong Kong University of Science and Technology (HKUST)**  
*Department of Computer Science and Engineering*  
Ph.D. in Computer Science and Engineering Sep. 2020 – Present  
◊ Advisor: Prof. Wei Wang (*expected to graduate in Fall 2025*)

**South China University of Technology (SCUT)**  
*School of Computer Science and Engineering*  
B.Eng. in Computer Science and Technology Sep. 2016 – Jul. 2020  
◊ Studied at All-English Innovation Class (GPA: 3.82/4); National Scholarship.

**PUBLICATIONS** \* *denotes co-first authors*  
**Refereed Papers in Conference Proceedings**  
[C5] Suyi Li\*, **Lingyun Yang**\*, Xiaoxiao Jiang, Hanfeng Lu, Dakai An, Zhipeng Di, Weiyi Lu, Jiawei Chen, Kan Liu, Yinghao Yu, Tao Lan, Guodong Yang, Lin Qu, Liping Zhang, Wei Wang, “Katz: Efficient Workflow Serving for Diffusion Models with Many Adapters,” in the *Proceedings of USENIX Annual Technical Conference (ATC ’25)*, Boston, MA, USA, July 2025. (*CCF-A, acceptance rate: 100/634=15.8%*)  
[C4] **Lingyun Yang**, Yongchen Wang, Yinghao Yu, Qizhen Weng, Jianbo Dong, Kan Liu, Chi Zhang, Yanyi Zi, Hao Li, Zechao Zhang, Nan Wang, Yu Dong, Menglei Zheng, Lanlan Xi, Xiaowei Lu, Liang Ye, Guodong Yang, Binzhang Fu, Tao Lan, Liping Zhang, Lin Qu, Wei Wang, “GPU-Disaggregated Serving for Deep Learning Recommendation Models at Scale,” in the *Proceedings of the 22nd USENIX Symposium on Networked Systems Design and Implementation (NSDI ’25)*, Philadelphia, PA, USA, April 2025. (*CCF-A, acceptance rate: 55/401=13.7%*)

- [C3] Qizhen Weng\*, **Lingyun Yang\***, Yinghao Yu, Wei Wang, Xiaochuan Tang, Guodong Yang, Liping Zhang, “Beware of Fragmentation: Scheduling GPU-Sharing Workloads with Fragmentation Gradient Descent,” in the *Proceedings of USENIX Annual Technical Conference (ATC ’23)*, Boston, MA, USA, July 2023. (*CCF-A*, acceptance rate: 65/353=18.4%)
- [C2] Yongkang Zhang, Yinghao Yu, Wei Wang, Qiukai Chen, Jie Wu, Zuowei Zhang, Jiang Zhong, Tianchen Ding, Qizhen Weng, **Lingyun Yang**, Cheng Wang, Jian He, Guodong Yang, and Liping Zhang, “Workload Management in Alibaba Clusters: The Good, the Bad, and the Ugly,” in the *Proceedings of ACM Symposium on Cloud Computing (SoCC ’22)*, San Francisco, CA, USA, November 2022. (*CCF-B*, acceptance rate: 38/155=24.5%)
- [C1] Luping Wang\*, **Lingyun Yang\***, Yinghao Yu, Wei Wang, Bo Li, Xianchao Sun, Jian He, and Liping Zhang, “Morphling: Fast, Near-Optimal Auto-Configuration for Cloud-Native Model Serving,” in the *Proceedings of ACM Symposium on Cloud Computing (SoCC ’21)*, Seattle, WA, USA, November 2021. (*CCF-B*, acceptance rate: 46/145=31.7%)  
**In submission / Preprint**
- [I1] Xiaoxiao Jiang\*, Suyi Li\*, Lingyun Yang, Tianyu Feng, Zhipeng Di, Weiyi Lu, Guoxuan Zhu, Xiu Lin, Kan Liu, Yinghao Yu, Tao Lan, Guodong Yang, Lin Qu, Liping Zhang, Wei Wang, “InstGenIE: Generative Image Editing Made Efficient with Mask-aware Caching and Scheduling,” *arXiv:2505.20600*, 2025.

|             |  |                     |
|-------------|--|---------------------|
| INTERNSHIP  | <b>Alibaba Group &amp; Alibaba Cloud</b>   | Hangzhou, China     |
|             | Research Intern  | Dec. 2020 – Present |
|             | ◊ Mentor: Dr. Yinghao Yu   |                     |
| 2024 – 2025 | <b><u>Efficient Text-to-Image Diffusion Model Serving with Add-on Modules</u></b>  |                     |
|             | ◊ Developed Katz, a system that efficiently generates high-quality images with stable diffusion models and many adapters (i.e., ControlNets and LoRAs).  |                     |
|             | ◊ Analyzed usage patterns of various adapters (ControlNets, LoRAs) based on <b>500k</b> request traces from production text-to-image services.   |                     |
|             | ◊ Incorporated several novel system designs, such as ControlNet-as-a-Service, bounded asynchronous LoRA loading, latent parallelism for CFG computation.   |                     |
|             | ◊ Achieved up to <b>7.8×</b> in serving latency and <b>1.6×</b> in throughput.   |                     |
|             | ◊ One paper was accepted by <b>ATC ’25</b> .   |                     |
| 2022 – 2024 | <b><u>GPU-Disaggregated Serving for Deep Learning Recommendation Models</u></b>  |                     |
|             | ◊ Proposed a GPU-disaggregated DLRM serving system to eliminate <i>resource mismatch</i> and meet <i>elastic</i> demand; leveraged RDMA network to offload computation on separate GPU and CPU nodes; resource-aware graph partitioning; topology-aware scheduling.      |                     |
|             | ◊ In daily scenarios (e.g., a crowded GPU cluster with > 90% allocation rate), reduced CPU fragments by <b>53%</b> and GPU fragments by <b>27%</b> . In the Double 11 Shopping Festival, saved up to <b>90%</b> of GPUs when loaning GPU servers from training clusters. |                     |
|             | ◊ One paper was accepted by <b>NSDI ’25</b> .  |                     |
| 2021 – 2023 | <b><u>Resource Fragmentation Analysis and Optimization for GPU-Sharing Clusters</u></b>  |                     |

- ◊ Formally quantified *statistical GPU resource fragments* and proposed the *fragmentation gradient descent* scheduling algorithm to reduce resource fragmentation. Our scheduling policy can significantly reduce *unallocated* GPUs by up to **49%** compared to state-of-the-art policies. [code] [trace]
- ◊ Developed ParaSet, a *best-effort* workload on Kubernetes that dynamically adjusts the number of instances and resource requirements based on the real-time resource availability in the cluster. It aims to fill resource fragments in the cluster and is integrated into KubeDL for internal use.
- ◊ One paper was accepted by ATC '23.

2021 – 2022

### **Large-Scale GPU Sharing and Overcommitment in Production**

- ◊ Enabled *large-scale GPU sharing and overcommitment* in production clusters, with *tens of thousands* of *shared* GPU containers running daily. Support the co-location of GPU tasks with different priorities (e.g., *latency-sensitive*, *spot*, *best-effort*).
- ◊ Specifically, I designed and implemented the *node-level* agent and the *cluster-level* controller. The agent periodically collects and reports resource usage metrics and *dynamically* allocates GPU resources to containers. The controller calculates potential resource overcommitment and provides scheduling guidance to the cluster scheduler.

2020 – 2021

### **Fast, Near-Optimal Auto-Configuration for Cloud-Native Model Serving**

- ◊ Developed Morphling, an auto-configuration framework for AI serving on Kubernetes; combined *meta-learning* and *bayesian optimization* to quickly find the *optimal* resource configuration (e.g., CPU cores, GPU timeshare, GPU memory, GPU type) and runtime parameters (e.g., batch size). [code]
- ◊ It was widely used in Alibaba for automated recommendation of container resource specifications; part of Alibaba's open-sourced KubeDL, a CNCF sandbox project.
- ◊ One paper was accepted by SoCC '21.

### **Microsoft Research Asia (MSRA)**

Beijing, China

Research Intern, *Innovation Engineering Group (IEG)*

Jul. 2019 – Jun. 2020

- ◊ Mentors: Lewei Lu & Chong Li
- ◊ Designed a novel pooling layer to enhance the robustness of image recognition models; built Neural Architecture Search (NAS) workflow for face recognition tasks; implemented various attention modules for face recognition models on CNTK.
- ◊ Star of Tomorrow Internship Award of Excellence.

### **AWARDS**

- ◊ USENIX NSDI 2025 Student Grant Mar. 2025
- ◊ Postgraduate Scholarship 2020 – 2024, HKUST
- ◊ Star of Tomorrow Internship Award of Excellence Jul. 2020, MSRA
- ◊ Merit Student & Excellent Student Cadre Nov. 2019, SCUT
- ◊ National Scholarship Oct. 2019, China
- ◊ Silver Medal, ICPC China Xi'an National Invitational Contest May 2019
- ◊ First Prize, 17th Guangdong Collegiate Programming Contest May 2019

|  |                 |
|--|-----------------|
| ◇ Silver Medal, 37Games Cup Programming Contest            | Apr. 2019       |
| ◇ Gold Medal, SCUT ACM Programming Contest                 | Apr. 2019       |
| ◇ Bronze Medal, ACM-ICPC Asia Xuzhou Regional Contest      | Oct. 2018       |
| ◇ Silver Medal, 1st Xiao Mi Collegiate Programming Contest | Sep. 2018       |
| ◇ Gold Medal, SCUT ACM Programming Contest                 | Apr. 2018       |
| ◇ The First Prize Scholarship                              | Nov. 2017, SCUT |
| ◇ Bronze Medal, ACM-ICPC Asia Xian Regional Contest        | Oct. 2017       |
| ◇ Gold Medal, 12th China Youth Robot Competition           | Jul. 2012       |
| ◇ Champion, RoboCup Youth Robot World Cup, China Division  | Mar. 2012       |

ACADEMIC  
SERVICES

**Artifact Evaluation Committee**

- ◇ ACM SIGCOMM (2024), IEEE HPCA (2024)
- ◇ ACM SOSP (2023), USENIX OSDI (2023), USENIX ATC (2023), MLSys (2023)

**External Conference Reviewer**

- ◇ IEEE INFOCOM (2022–2025), IEEE ICDCS (2023, 2025)
- ◇ ACM APSys (2021), IEEE MSN (2021), EAI Qshine (2020)

**External Journal Reviewer**

- ◇ IEEE Transactions on Cloud Computing (2025)

**Student Volunteer**

- ◇ ACM APNet (2023), IEEE ICMLC & ICWAPR (2018)

TEACHING  
ACTIVITIES

**Hong Kong University of Science and Technology**

*Teaching Assistant, Department of Computer Science and Engineering*

- ◇ CSIT6000O: Advanced Cloud Computing (Spring 2022, Spring 2023)
- ◇ COMP4651: Cloud Computing and Big Data Systems (Spring/Fall 2021, Spring 2024)
- ◇ COMP3511: Operating Systems (Fall 2023)

UNDERGRADUATE  
EXPERIENCE

**ACM-ICPC Competition Group**

SCUT

*Group Member & Team Leader*

2016 – 2019

- ◇ Coach: Prof. Chuhua Xian
- ◇ Major domains: Dynamic Programming, Number Theory, Data Structure, etc.

**Machine Learning & Cybernetics Research Group**

SCUT

*Undergraduate Research Assistant*

2017 – 2019

- ◇ Advisor: Prof. Patrick Chan

- ◇ Projects: Fundus Stitching, Tableware Recognition, Neural Network Visualization.

**Tencent Innovation Club**

SCUT, CSE

*Vice Chairman*

2018 – 2019

- ◇ Led the *largest* student club in SCUT CSE, sponsored by Tencent.

**ByteDance Summer Camp**

Beijing, China

*Camper, Algorithm track*

Aug. 2019

- ◇ Mentor: Dr. Yibo Zhu

◇ Totally 150 participants selected from more than 6k candidates ( $< 2.5\%$ ).

SKILLS

Programming Languages: Golang, C++, Python, Javascript

Toolkits: Kubernetes, Docker, Grafana, Git,  $\LaTeX$ , SQL, Markdown

Languages: English (fluent), Mandarin (Native speaker), Cantonese (Intermediate)

MISCELLANEOUS

Play basketball & badminton & squash, workout at the gym, foodie.

MISC

My paper reading notes are available at <https://paper.lingyunyang.com/>.