

# Lingyun Yang (杨凌云)

Updated July 23, 2024

**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/>

## PROFILE

Ph.D. Candidate  
Department of Computer Science and Engineering  
Hong Kong University of Science and Technology  
Clear Water Bay, Kowloon, Hong Kong

## RESEARCH INTERESTS

I am broadly interested in resource management for large-scale data centers, especially in improving resource efficiency for large-scale AI/GPU clusters and optimizing model serving systems.

## EDUCATION

### **Hong Kong University of Science and Technology (HKUST)**

*Department of Computer Science and Engineering*

Ph.D. in Computer Science and Engineering 2020 – Present

◊ Advisor: Prof. Wei Wang

### **South China University of Technology (SCUT)**

*School of Computer Science and Engineering*

B.Eng. in Computer Science and Technology 2016 – 2020

◊ Studied at All-English Innovation Class (GPA: 3.82/4)

## PROFESSIONAL EXPERIENCE

### **Alibaba Group**

Hangzhou, China

*Research Intern, Cluster Management Group*

Dec. 2020 – Present

◊ Mentor: Dr. Yinghao Yu

### **Mitigate GPU Resource Fragmentation**

◊ **TBD.** GPU disaggregation for recommendation services.

◊ [C3] Formally defined *statistical GPU resource fragments* and proposed the *fragmentation gradient descent* scheduling algorithm to reduce resource fragmentation. Large-scale trace evaluations show that our scheduling policy can significantly improve GPU allocation rate by 3% compared to state-of-the-art policies.

◊ Developed *ParaSet*, a best-effort workload on Kubernetes that can dynamically adjust 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. It is integrated into KubeDL for internal use.

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

◊ Enabled *large-scale GPU sharing* in production clusters, with over 10k shared GPU containers running on a daily basis. Support the co-location of GPU tasks with different priorities (e.g., *latency-sensitive*, *best-effort*).

◊ This was a multi-team collaborative project. I was responsible for the design and implementation of the single-node agent and the centralized controller. The agent periodically collects and reports resource usage metrics, as well as dynamically allocates GPU resources to containers. The controller calculates potential overcommitment of resources and provides scheduling guidance to the scheduler.

### **Efficient Diffusion Model Serving System**

◊ TBD.

### **Auto-Configuration for AI Serving Service**

◊ [C1] Developed Morphling, an open-source auto-configuration framework for AI serving on Kubernetes. It combines meta-learning and bayesian optimization to quickly find the optimal configuration. Internally, it was widely used for automated recommendation of container resource specifications. [[code](#)]

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

Beijing, China

*Research Intern*, Innovation Engineering Group (IEG)

Jul. 2019 – Jun. 2020

◊ Research on model robustness, face recognition, attention mechanisms, knowledge distillation, and neural architecture search.

### **PUBLICATIONS**

\* *denotes co-first authors*

[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.

[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.

[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.

*Preprint*

[P1] Suyi Li, **Lingyun Yang\***, Xiaoxiao Jiang, Hanfeng Lu, Zhipeng Di, Weiyi Lu, Jiawei Chen, Kan Liu, Yinghao Yu, Tao Lan, Guodong Yang, Lin Qu, Liping Zhang, Wei Wang, “SwiftDiffusion: Efficient Diffusion Model Serving with Add-on Modules,” *arXiv preprint arXiv:2407.02031*, 2024.

### **HONORS AND SCHOLARSHIPS**

◊ Postgraduate Scholarship	2020 – Present, 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 Xian National Invitational Contest	May 2019
◊ First Prize, 17th Guangdong Collegiate Programming Contest	May 2019

	<ul style="list-style-type: none"> <li>◊ Silver Medal, 37Games Cup Programming Contest Apr. 2019</li> <li>◊ Gold Medal, SCUT ACM Programming Contest Apr. 2019</li> <li>◊ Bronze Medal, ACM-ICPC Asia Xuzhou Regional Contest Oct. 2018</li> <li>◊ Silver Medal, 1st Xiao Mi Collegiate Programming Contest Sept. 2018</li> <li>◊ Gold Medal, SCUT ACM Programming Contest Apr. 2018</li> <li>◊ The First Prize Scholarship Nov. 2017, SCUT</li> <li>◊ Bronze Medal, ACM-ICPC Asia Xian Regional Contest Oct. 2017</li> <li>◊ Gold Medal, 12th China Youth Robot Competition Jul. 2012</li> <li>◊ Champion, RoboCup Youth Robot World Cup, China Division Mar. 2012</li> </ul>
ACADEMIC SERVICES	<p><b>Artifact Evaluation Committee</b></p> <ul style="list-style-type: none"> <li>◊ SIGCOMM (2024), HPCA (2024)</li> <li>◊ SOSP (2023), OSDI (2023), ATC (2023), MLSys (2023)</li> </ul> <p><b>External Reviewer</b></p> <ul style="list-style-type: none"> <li>◊ INFOCOM (2022, 2023, 2024)</li> <li>◊ ICDCS (2023), APSys (2021), MSN (2021), Qshine (2020)</li> </ul> <p><b>Student Helper</b></p> <ul style="list-style-type: none"> <li>◊ APNet (2023), ICMLC &amp; ICWAPR (2018)</li> </ul>
TEACHING ACTIVITIES	<p><b>Hong Kong University of Science and Technology</b></p> <p><i>Teaching Assistant, Department of Computer Science and Engineering</i></p> <ul style="list-style-type: none"> <li>◊ CSIT6000O: Advanced Cloud Computing (Spring 2022, Spring 2023)</li> <li>◊ COMP4651: Cloud Computing and Big Data Systems (Spring 2021, Fall 2021, Spring 2024)</li> <li>◊ COMP3511: Operating Systems (Fall 2023)</li> </ul>
OTHER EXPERIENCE	<p><b>ACM-ICPC Competition Group</b></p> <p><i>Group Member &amp; Team Leader</i> 2016 – 2019</p> <ul style="list-style-type: none"> <li>◊ Coach: Prof. Chuhua Xian</li> <li>◊ Major domains: Dynamic Programming, Number Theory, Data Structure, etc.</li> </ul> <p><b>Machine Learning &amp; Cybernetics Research Group</b></p> <p><i>Undergraduate Research Assistant</i> 2017 – 2019</p> <ul style="list-style-type: none"> <li>◊ Advisor: Prof. Patrick Chan</li> <li>◊ Projects: Fundus Stitching, Tableware Recognition, and NN Visualization.</li> </ul> <p><b>Tencent Innovation Club</b></p> <p><i>Vice Chairman</i> 2018 – 2019</p> <ul style="list-style-type: none"> <li>◊ Led the largest student club in SCUT CSE, sponsored by Tencent.</li> </ul> <p><b>ByteDance Summer Camp</b> Beijing, China</p> <p><i>Camper, Algorithm track</i> Aug. 2019</p> <ul style="list-style-type: none"> <li>◊ Mentor: Dr. Yibo Zhu</li> <li>◊ Totally 150 participants selected from more than 6k candidates.</li> </ul>
SKILLS	<p>Programming Languages: Golang, C++, Python</p> <p>Toolkits: Kubernetes, Git, <math>\text{\LaTeX}</math>, Linux Shell, Qt, MySQL, Markdown</p>

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

MISCELLANEOUS    Play basketball & badminton & squash, workout at the gym, foodie.  
My paper reading notes are available at <https://paper.lingyunyang.com/>.