

# Lingyun Yang (杨凌云)

Updated August 11, 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 have a broad interest in resource management for large-scale data centers. Specifically, my research focuses on: (a) improving <i>resource efficiency</i> for AI/GPU clusters; (b) building <i>efficient</i> and <i>low-cost</i> AI model serving systems.
EDUCATION	<b>Hong Kong University of Science and Technology (HKUST)</b> <i>Department of Computer Science and Engineering</i> Ph.D. in Computer Science and Engineering 2020 – Present ◊ Advisor: Prof. Wei Wang ( <i>expected to graduate in Fall 2025</i> )  <b>South China University of Technology (SCUT)</b> <i>School of Computer Science and Engineering</i> B.Eng. in Computer Science and Technology 2016 – 2020 ◊ Studied at All-English Innovation Class (GPA: 3.82/4)
PROFESSIONAL EXPERIENCE	<b>Alibaba Group</b> Hangzhou, China <i>Research Intern, Cluster Management Group</i> Dec. 2020 – Present ◊ Mentor: Dr. Yinghao Yu <b>Resource Management for AI/GPU Clusters</b> <u>Mitigate GPU Resource Fragmentation</u> ◊ [ <b>Under Review</b> ] Proposed the GPU-disaggregated DLRM serving system to eliminate <i>resource mismatch</i> and meet <i>elastic</i> demand. By leveraging RDMA network to compute the computation graph on GPU and CPU nodes separately, it reduced CPU fragments by 53% and GPU fragments by 27%. Saved up to 90% of GPUs when loaning GPU servers from training clusters, during the seasonal traffic peaks (e.g., Double 11 Shopping Festival). ◊ [ <b>ATC 2023</b> ] Formally quantified <i>statistical GPU resource fragments</i> and proposed the <i>fragmentation gradient descent</i> scheduling algorithm to reduce resource fragmentation. Our scheduling policy can significantly reduce <i>unallocated</i> GPUs by up to 49% compared to state-of-the-art policies. [ <a href="#">code</a> ]

- ◊ 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 and 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*).
- ◊ Specifically, I designed and implemented the node-level agent and the cluster-level controller. The agent periodically collects and reports resource usage metrics, as well as dynamically allocates GPU resources to different containers. The controller calculates potential resource overcommitment and provides scheduling guidance to the cluster scheduler.

### **Efficient and Low-cost AI Model Serving Systems**

#### Efficient Diffusion Model Serving with Add-on Modules

- ◊ [**Under Review**] Developed SwiftDiffusion, a system that efficiently generates high-quality images with stable diffusion models and add-on modules (i.e., ControlNets and LoRAs). Incorporated several novel designs, including ControlNet-as-a-Service, asynchronous LoRA loading, and kernel optimization. Achieved up to 5× in latency and 2× in throughput without sacrificing image quality.

#### Auto-Configuration for AI Serving Service

- ◊ [**SoCC 2021**] Developed Morphling, an open-source auto-configuration framework for AI serving on Kubernetes. Combined *meta-learning* and *bayesian optimization* to quickly find the *optimal* configuration. It was widely used in Alibaba 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*

- ◊ 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.

- ◊ **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,” *under review*.

- ◇ 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.
- ◇ 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.
- ◇ 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.

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
	◇ 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	Sept. 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	<b>Artifact Evaluation Committee</b>	
	◇ SIGCOMM (2024), HPCA (2024)	
	◇ SOSP (2023), OSDI (2023), ATC (2023), MLSys (2023)	
	<b>External Reviewer</b>	
	◇ INFOCOM (2022, 2023, 2024)	
	◇ ICDCS (2023), APSys (2021), MSN (2021), Qshine (2020)	
	<b>Student Helper</b>	
	◇ APNet (2023), ICMLC & ICWAPR (2018)	

TEACHING ACTIVITIES	<b>Hong Kong University of Science and Technology</b>	
	<i>Teaching Assistant, Department of Computer Science and Engineering</i>	
	◇ CSIT6000O: Advanced Cloud Computing (Spring 2022, Spring 2023)	
	◇ COMP4651: Cloud Computing and Big Data Systems (Spring 2021, Fall 2021, Spring 2024)	

◇ COMP3511: Operating Systems (Fall 2023)

OTHER EXPERIENCE	<b>ACM-ICPC Competition Group</b>	SCUT
	<i>Group Member &amp; Team Leader</i>	2016 – 2019
	◇ Coach: Prof. Chuhua Xian	
	◇ Major domains: Dynamic Programming, Number Theory, Data Structure, etc.	
	<b>Machine Learning &amp; Cybernetics Research Group</b>	SCUT
	<i>Undergraduate Research Assistant</i>	2017 – 2019
	◇ Advisor: Prof. Patrick Chan	
	◇ Projects: Fundus Stitching, Tableware Recognition, and NN Visualization.	
	<b>Tencent Innovation Club</b>	SCUT, CSE
	<i>Vice Chairman</i>	2018 – 2019
SKILLS	◇ Led the <i>largest</i> student club in SCUT CSE, sponsored by Tencent.	
	<b>ByteDance Summer Camp</b>	Beijing, China
	<i>Camper</i> , Algorithm track	Aug. 2019
MISCELLANEOUS	◇ Mentor: Dr. Yibo Zhu	
	◇ Totally 150 participants selected from more than 6k candidates.	
	Programming Languages: Golang, C++, Python, Javascript Toolkits: Kubernetes, Docker, Grafana, Git, <del>LaTeX</del> $\text{\LaTeX}$ , SQL, Markdown Languages: English (fluent), Mandarin (Native speaker), Cantonese (Intermediate)	
MISCELLANEOUS	Play basketball & badminton & squash, workout at the gym, foodie.	
	My paper reading notes are available at <a href="https://paper.lingyunyang.com/">https://paper.lingyunyang.com/</a> .	