

# Lingyun Yang (杨凌云)

Updated May 31, 2023

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

**Phone:** (+852) 9719-0933

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

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

**LinkedIn:** [lingyun-yang-b2881b18a](https://www.linkedin.com/in/lingyun-yang-b2881b18a)

**Web:** <https://www.cse.ust.hk/~lyangbk/>

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 production clusters, especially in improving resource efficiency for AI/GPU clusters and optimizing system performance bottlenecks using machine learning methods.	
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	
	<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>Reduce GPU Resource Fragmentation</b>	
	◊ [C3] Formally defined <i>GPU resource fragments</i> and proposed the <i>fragmentation gradient descent</i> algorithm to reduce resource fragmentation during scheduling. 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 <i>ParaSet</i> , 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.	
	<b>Large-Scale GPU Sharing in Production</b>	
	◊ Enabled <i>large-scale GPU sharing</i> in production clusters, with over 4000 shared GPU containers running on a daily basis. Support the co-location of GPU tasks with different priorities (e.g., <i>latency-sensitive</i> , <i>best-effort</i> ).	

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

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

◇ [C1] Co-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 is 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 network search.

#### **PUBLICATIONS**

[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,” to appear in the *Proceedings of USENIX Annual Technical Conference (ATC ’23)*, Boston, MA, USA, July 2023. (\*Co-first authors in alphabetical order)

[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. (\*Co-first authors in alphabetical order)

#### **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

	<ul style="list-style-type: none"> <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>◇ OSDI (2023), ATC (2023), MLSys (2023)</li> </ul> <p><b>External Reviewer</b></p> <ul style="list-style-type: none"> <li>◇ INFOCOM (2022, 2023), ICDCS (2023), APSys (2021), MSN (2021), Qshine (2020)</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)</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> <p><b>ICMLC &amp; ICWAPR</b> Chengdu, China</p> <p><i>Student helper</i> Jul. 2018</p>
SKILLS	<p>Programming Languages: Golang, C++, Python</p> <p>Toolkits: Kubernetes, Git, <math>\LaTeX</math>, Linux Shell, Qt, MySQL, Markdown</p> <p>Languages: English (fluent), Mandarin (Native speaker)</p>
MISCELLANEOUS	<p>Play basketball &amp; badminton, workout at the gym, food lover.</p>