

# Zhisheng YE

Ph.D. Student

EECS, Peking University  
1336, 1st Science Building, No.5 Yiheyuan Lu  
Haidian District, Beijing, 100871, P. R. China  
✉ yezhisheng@pku.edu.cn  
🏠 <https://yeshisheng.me> · 🌐 yzs981130 · 📄 yzs981130

## RESEARCH INTERESTS

I am interested in Distributed Systems, Systems for Machine Learning and Resource Management. My current focuses include:

- Scheduling systems for Deep Learning training workloads.
- Elastic & flexible resource management on heterogeneous clusters.
- Co-optimizing deep learning framework with scheduling.

## EDUCATION

**Peking University** Beijing, China  
Pursuing Ph.D. in Computer Science *Sep 2019 – Present*  
Related Courses: Distributed Systems / Introduction to Parallel Computing / Dataflow Algorithms  
Advisor: Prof. Yingwei Luo & Prof. Xiaolin Wang

**Peking University** Beijing, China  
Bachelor of Science in Computer Science & Technology *Sep 2015 – Jun 2019*  
Related Courses: Computer Network (Honor track) / Computer Network Practice (Honor track) / Operating System (Honor track) / Introduction to Computer System

## PROFESSIONAL EXPERIENCE

**Peng Cheng Laboratory** Shenzhen, China  
**Research Intern** *July 2018 – Present*

- Contributed to development of OpenI-Octopus, an open-sourced scheduler for deep learning training workloads based on Kubernetes.
- Safe GPU sharing mechanisms for multiple workloads
  - Intercepted job usage of GPU by deep learning framework analysis and CUDA driver / library wrapping to securely limit job usage of GPU memory.
  - Open-sourced the mechanisms of safe GPU sharing under memory limitation.
- GPU sharing techniques on Kubernetes
  - Implemented resource allocation and usage techniques of vGPU on Kubernetes as a device plugin. The division of GPUs can be pre-defined or set dynamically.
  - Implemented a monitoring system of vGPU runtime statistics based on Prometheus.
- Monitoring & logging systems on Kubernetes
  - Implemented job lifecycle monitoring and event logging system on Kubernetes.
  - Collected and analyzed traces of deep learning workloads in PCL cluster.
  - Modeled the characteristics of deep learning workloads based on the trace analysis of fine-grained GPU resource usage.
- Next generation GPU sharing mechanisms
  - Supported high-efficient pausing and resuming of deep learning workloads, which is imperceptible to users.
  - Current working on workloads migration and further combination with scheduling algorithms.
- Developed a flow-based scheduler on Kubernetes, supporting batch scheduling of deep learning workloads.

**Sensetime Research** Beijing, China  
**Research Intern** *Sept 2019 – Present*

- Advisor: Peng Sun from Sensetime Research. I am also co-advised by Tianwei Zhang from S-Lab, Nanyang Technological University.
- Creativity projects
  - Developed and drafted Astraea, a fair scheduler for deep learning training jobs, as first author. Currently under submission.

- Authored 4 patents related to scheduling as first author.
- Communication optimization of GPUs across pods
  - Modified the logic of kubelet, communication framework and application framework to use GPU, to achieve cross-pod P2P communication
  - Achieved near bare-metal performance; 35% improvement over original.
- Designed and implemented a production scheduler on Kubernetes, introducing GPU topology-aware scheduling.
- Cluster trace analysis and characteristics of deep learning workloads, aiming on pending situation and daily / monthly trends of job behaviours.
- Implemented and designed several key components of infrastructure in Sensetime, including container images for business lines, a Golang module server and security enhancements in Kubernetes.

**Student Super Computing Team of Peking University**  
**Team member**

Beijing, China  
 Sept 2018 – June 2019

- Advisor: Prof. Xiaolin Wang & Prof. Yun (Eric) Liang
- Participated in analyzing, compiling, profiling, optimizing, and improving parallelizability of general HPC tasks.
- Optimized the compilation and operation of a specific supercomputing task: the gene assembly problem wtdbg2 in ASC19 Student Supercomputer Challenge; Participated and led the group competition challenge in the final round.

**SKILLS**

**Programming** C/C++, Golang, Python, Bash, Javascript  
**Software & Tools** Kubernetes, CUDA  
**Languages** English, Mandarin

**PUBLICATIONS**

1. **Z. Ye**, P. Sun, W. Gao, T. Zhang, X. Wang, S. Yan, and Y. Luo, “Astraea: A fair deep learning scheduler for multi-tenant gpu clusters,” *IEEE Transactions on Parallel and Distributed Systems (TPDS)*, 2021.
2. W. Gao, **Z. Ye**, P. Sun, Y. Wen, and T. Zhang, “Chronus: A novel deadline-aware scheduler for deep learning training jobs,” in *Proceedings of the ACM Symposium on Cloud Computing (SoCC '21)*, Association for Computing Machinery, 2021.

**TEACHING**

- Teaching Assistant: Introduction to Computing (Peking University) 2019 Fall, 2020 Fall
- Teaching Assistant: Design and Analysis of Algorithms (Peking University) 2020 Spring, 2021 Spring
- Teaching Assistant: Introduction to Virtualization and Storage Systems (Peking University) 2019 Fall

**AWARDS**

- First Price (Team), ASC Student Supercomputer Challenge 2019
- Outstanding Winner, Microsoft Student Club Practice Space 2018