

Chenxia HAN

+ (86) 176-1131-8773 | cxhan@link.cuhk.edu.hk | cxhan.com

EDUCATION

The Chinese University of Hong Kong <i>M.Phil. in Computer Science and Engineering</i>	Hong Kong 2020 - Nov. 2024
Wuhan University <i>B.Eng. in Computer Science and Technology</i>	Wuhan 2014 - 2018

RESEARCH

My previous work has primarily focused on three key areas: (1) developing efficient training frameworks, (2) accelerating model inference for video processing, and (3) optimizing GPU resource management.

SELECTED PUBLICATIONS

Scalable Complex Event Processing on Video Streams Chenxia Han , Chaokun Chang, Srijan Srivastava, et al.	<i>SIGMOD 2025</i>
SGDRC – Software-Defined Dynamic Resource Control for Concurrent DNN Inference on NVIDIA GPUs Yongkang Zhang, Haoxuan Yu, Chenxia Han , et al.	<i>PPoPP 2025</i>
Top-K Deep Video Analytics: A Probabilistic Approach Ziliang Lai, Chenxia Han , Chris Liu, et al.	<i>SIGMOD 2021</i>
SimpleDet: A Simple and Versatile Distributed Framework for Object Detection and Instance Recognition Yuntao Chen, Chenxia Han , Yanghao Li, et al.	<i>JMLR 2019</i>

WORKING EXPERIENCES

Research Assistant <i>The Chinese University of Hong Kong</i>	Hong Kong Aug. 2019 - Nov. 2024
<ul style="list-style-type: none">Developed Bobsled, an efficient system for video stream pattern matching. The solution combines algorithmic innovations (<u>a draft model to reduce target model inference</u>) with system optimizations (<u>code generation, caching, and batching</u>) to minimize overhead. Achieved speedups of $2.4\times$ to $11.6\times$ without accuracy degradation.Created SGDRC, a dynamic GPU resource controller that <u>optimizes VRAM allocation</u> by reverse-engineering black-box channel mappings. This novel approach <u>eliminates resource contention</u> during concurrent model inference, delivering throughput improvements up to $1.47\times$ compared to existing GPU sharing solutions.Designed Everest, a high-performance system for video top-K query processing. Our solution combines a draft model with a probabilistic algorithm to achieve $14.3\times$ to $20.6\times$ speedups over traditional approaches.	
Research Intern <i>TuSimple</i>	Beijing May 2018 - June 2019
<ul style="list-style-type: none">Developed SimpleDet (3.1k ☆), a high-performance object detection framework built on MXNet. The system implements cutting-edge optimizations including <u>mixed-precision training</u>, <u>cross-GPU synchronized batch normalization</u>, <u>gradient checkpointing</u>, and <u>operator fusion</u>, delivering $1.7\times$ faster Mask R-CNN training throughput in FP16 compared to existing frameworks.	

HONORS AND AWARDS

Best Demo Paper Runner Up , SIGMOD	2022
Postgraduate Scholarship , CUHK	2020-2022
Gold Medal , Google AI Open Images Object Detection Track	2018
Silver Medal , China Collegiate Programming Contest	2015

SKILLS

Programming Languages	C/C++, CUDA, Python, Golang, Java, Scala
Frameworks	PyTorch, MXNet, TVM, TensorRT, Horovod, MPI, Ray, Flink