

# KEHAN QI

Ph.D. Student | Eligible for CPT/OPT  
kehan.qi@stonybrook.edu | (934)263-2748

## EDUCATION

### Stony Brook University

Ph.D. Student

Biomedical Informatics

Aug 2024 - Aug 2028 (expected)

### University of Chinese Academy of Sciences

Master of Engineering

Computer Technology

Sept 2018 - June 2021

### Zhejiang University

Bachelor of Engineering

Measurement Control Technology and Instruments

Sept 2013 - July 2017

## WORK EXPERIENCE

### Stony Brook University

Research Assistant

Stony Brook, NY, US

May 2025 - present

- **Flow-based MRI Reconstruction:** Design and implement conditional probability flow for MRI reconstruction using flow matching. Paper submitted to ICLR 2026.
- **VLM-based Pathology Image Processing:** LLaVA with learnable tokens to compress giant number of tokens for Whole-Slide Image-guided text generation. [Preprint PDF]
- **LLM-based Paper Review BOT (in-progress):** LLM & RAG & open dataset for paper review; end-to-end BOT deployed on HF and AWS. Progress: Implemented MVP; building up V1 and dataset. [PRD][System Design]

### Stori

Data Engineer (MLE focus)

Full-time Employee, Hangzhou, China

Apr 2023 - July 2024

- **Real-time Data Infrastructure:** Architected and built real-time data pipelines using Flink, Kinesis, Lambda, DynamoDB, and Elasticsearch, enabling consistent pipelines for ML model invocation, data monitoring, and downstream query API integration.
- **Low-latency ML Inference System for Risk Control:** Owned end-to-end design, implementation, and deployment of a real-time ML inference system for transaction-level risk control. Used AWS DMS + Flink + Kinesis + Lambda + SageMaker for cross-account model invocation with 100ms average latency.
- **Team Leadership and Standards:** Established internal coding and deployment standards, CI/CD pipeline, and AWS CDK infrastructure templates. Mentored two junior engineers and led weekly sprint planning and code reviews.

### Amazon

Software Development Engineer (Applied ML System)

Full-time Employee, Beijing, China

Aug 2021 - Feb 2023

- **Applied ML System Engineering:** Designed and implemented an automated pipeline for scheduled product classification updates using ML models deployed on AWS SageMaker. Integrated Lambda, SNS, S3, and DynamoDB to support scalable, production-level ML inference and ingestion with tens of millions of products.
- **Impact Analysis via Distributed Processing:** Built large-scale data process pipelines to evaluate financial impacts of updated classification models using Spark. Analyzed ~10B records to compute fee deltas pre- and post-deployment across multiple dimensions (product, seller, category). Applied Spark job optimization to reduce runtime to within 20 minutes.
- **Future Fee Prediction System:** Developed inference-based fee projection system utilizing classification results. Performed batch processing on 1.5B+ records with AWS Glue and Redshift, and optimized TPS throttling to support SageMaker-based fee computation. Enabled daily updates within a 24-hour SLA.

### Tencent

Research Intern

Intern, Shenzhen, China

June 2020 - Sept 2020

- **Registered medical image quality analysis:** a) Detect landmarks from registered CT images. b) Train a neural network to predict registered image quality score, with landmarks and registered image as input. c) A Chinese patent produced.

## SELECTED PAPERS

- **Kehan Qi**, Saumya Gupta, Qingqiao Hu, Weimin Lyu, and Chao Chen\*. "Unrolled Networks Are Conditional Probability Flow ODEs in MRI Reconstruction". (ICLR 2026 under review)
- Weimin Lyu, Qingqiao Hu, **Kehan Qi**, Zhan Shi, Wentao Huang, Saumya Gupta, and Chao Chen\*. "Efficient Whole Slide Pathology VQA via Token Compression." arXiv preprint arXiv:2507.14497 (2025).
- **Kehan Qi**, Hao Yang, Cheng Li, Zaiyi Liu, Meiyun Wang, Qiegen Liu, and Shanshan Wang\*. "X-Net: Brain Stroke Lesion Segmentation Based on Depthwise Separable Convolution and Long-range Dependencies". MICCAI 2019.
- Hao Yang, Weijian Huang, **Kehan Qi**, Cheng Li, Xinfeng Liu, Meiyun Wang, Hairong Zheng, and Shanshan Wang\*. "CLCI-Net: Cross-Level Fusion and Context Inference Networks for Lesion Segmentation of Chronic Stroke". MICCAI 2019.

## SKILLS

- **Data Processing Techniques:** Spark, Flink, Hive, MySQL, No-SQL
- **Amazon Web Service (AWS) Skills:** Glue, EMR, Lambda Function, SQS, Managed Service for Apache Flink, API Gateway, VPC, DMS, S3, SageMaker
- **Deep Learning Techniques:** Diffusion Model, In-context Learning, Visual Language Model