

KEHAN QI

Ph.D. Student | kehan.qi@stonybrook.edu | (934)263-2748
andrewsher.github.io

EDUCATION

Stony Brook University

Ph.D. Student in Biomedical Informatics

Stony Brook, NY, US

Aug 2024 - Aug 2028 (expected)

University of Chinese Academy of Sciences

M. Eng. in Computer Technology

Shenzhen, China

Sep 2018 - Jun 2021

Zhejiang University

B. Eng. in Measurement Control Technology and Instruments

Hangzhou, China

Sep 2013 - Jul 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.
- **VLM-based Pathology Image Processing:** LLaVA with learnable tokens to compress giant number of tokens for Whole-Slide Image-guided text generation. [PDF]
- **LLM-based Paper Review BOT:** 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)

Hangzhou, China

Apr 2023 - Jul 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)

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

Shenzhen, China

Jun 2020 – Sep 2020

- **Medical Image Processing Research:** Conducted research on CT image processing methods for registration and quality assessment.
- **Registration Quality Assessment:** Developed a landmark-based neural network for evaluating registration quality of medical images.

SELECTED PAPERS

- Qingqiao Hu, Weimin Lyu, Meilong Xu, **Kehan Qi**, Xiaoling Hu, and Chao Chen*. "LoC-Path: Learning to Compress for Pathology Multimodal Large Language Models." (Under Review)
- **Kehan Qi**, Saumya Gupta, Qingqiao Hu, Weimin Lyu, and Chao Chen*. "Unrolled Networks Are Conditional Probability Flow ODEs in MRI Reconstruction". (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, API Gateway, VPC, DMS, S3, SageMaker
- **Deep Learning Techniques:** Diffusion Model, In-context Learning, Visual Language Model