# QIRUN DAI

(+1) 530-709-2773 ⋄ qirundai@illinois.edu/daiqirun19@gmail.com ⋄ Website

## EDUCATION

**Fudan University, Shanghai**                                            *Sept. 2021 - Present*
*B.Eng. in Artificial Intelligence (Honor Class, Data Science Track)*     **GPA:** 3.89/4.00; **Rank:** 1/25
**Course Highlights:** Mathematical Analysis I/II/III (A), Advanced Linear Algebra (A), Methods of Optimization (A), Stochastic Processes (A), Data Fusion and Assimilation (A), Set Theory and Graph Theory (A), Data Structures and Algorithm Design (A), Artificial Intelligence (A), Natural Language Processing (A)

**University of California, Davis**                                       *Sept. 2023 - Dec. 2023*
*Exchange Student in Computer Science*                                    **GPA:** 4.00/4.00
**Course Highlights:** Operating Systems (A), Numerical Linear Algebra (Graduate, A), Advanced Statistical Learning (Graduate, A+)

## RESEARCH INTERESTS

My research interests span **natural language processing** and **machine learning**. These days, I have an intense interest in scientifically advancing the capabilities, efficiency and trustworthiness of large language models (LLMs) from a **data-centric** perspective, with a specific focus on these topics:

- Data Curation and Synthetic Data.

- Data-efficient Learning Algorithms.

- Data Attribution for Responsible and Interpretable Models.

- Datasets for Trustworthy and Efficient Evaluation.

## PUBLICATIONS

**Improving Influence-based Instruction Tuning Data Selection for Balanced Learning of Diverse Capabilities**
**Qirun Dai**, Dylan Zhang, Jiaqi W. Ma, Hao Peng.
*Under Review at ICLR 2025*

**Demonstration Distillation for Efficient In-Context Learning**
Tong Chen, **Qirun Dai**, Zhijie Deng, Dequan Wang.
*Submitted to ICLR 2024*

## RESEARCH EXPERIENCE

**University of Illinois Urbana-Champaign**                               *Apr. 2024 - Present*
Research Intern, Siebel School of Computing and Data Science

- **Improving Influence-based Data Selection for Multi-task Instruction Tuning**
  *Advisor: Prof. Hao Peng and Prof. Jiaqi W. Ma*                        *Apr. 2024 - Oct. 2024*

  – When instruction tuning LLMs for learning multiple diverse tasks, we identified the poor performance of data selection methods built upon gradient-based influence estimation techniques, and attributed this problem to an inherent bias in cross-task influence.

  – We then proposed BIDS, a **B**alanced and **I**nfluential **D**ata **S**election algorithm that addresses this problem through instance-level normalization and iterative selection.

– When training on UltraInteract, a SOTA high-quality dataset designed to enhance diverse reasoning capabilities, we showed that a 15% subset selected by BIDS can outperform full-dataset training in terms of the overall performance on 7 benchmarks spanning coding, math, STEM, logical reasoning and instruction following. We provided in-depth analyses on what might be the good properties of a balanced set of influential data.

– This work resulted in a first-authored paper currently under review at ICLR 2025.

### Shanghai Jiao Tong University
*Jan. 2023 - Mar. 2024*
Research Intern, Qing Yuan Research Institute

- **Demonstration Selection for Knowledge- & Reasoning-Intensive In-Context Learning**
  *Advisor: Prof. Dequan Wang and Prof. Zhijie Deng*      *Nov. 2023 - Mar. 2024*

  – Proposed an ICL demonstration selection method built on sparse retrieval techniques such as BM25, targeting knowledge- and reasoning-intensive QA tasks where pretrained dense embeddings can be easily confused due to a lack of strong reasoning or domain-specific knowledge.

  – Achieved more than 4% accuracy improvement in multiple challenging domains including medicine and college mathematics, which were rarely explored by previous ICL research.

  – Explored what makes a good demonstration for improving inference-time learning efficiency.

- **Demonstration Distillation for Efficient In-Context Learning**
  *Advisor: Prof. Dequan Wang and Prof. Zhijie Deng*      *July 2023 - Sept. 2023*

  – To optimize context efficiency for LLM in-context learning (ICL), we developed DGS, an agentic framework that iteratively compresses demonstrations with three LLM agents: **D**istillist, **G**eneralist and **S**pecialist.

  – DGS achieved up to a 4.3x compression ratio and a 5% accuracy improvement on various QA tasks, showing that overly concise demonstrations under human perception can still be informative for LLMs to learn at inference time, and elicit even stronger reasoning capabilities.

  – This work resulted in a second-authored paper submitted to ICLR 2024.

- **The Forward-Forward Algorithm: Some In-Depth Investigations**
  *Advisor: Prof. Dequan Wang*      *Jan. 2023 - May 2023*

  – Generalized the Forward-Forward learning Algorithm (FFA) proposed by Geoffrey Hinton to modern CNNs and Vision Transformers.

  – Inspired by classical contrastive learning paradigms, explored different methods to generate positive and negative training data, and how they boost the learning efficiency of FFA.

  – To tackle the inherent training instability of FFA, conducted ablation studies on various factors including residual connection, normalization methods, loss function design and learnability of threshold values, resulting in an FFA training recipe especially for modern vision models.

## TEACHING EXPERIENCE

### Introduction to Computer Systems (DATA130025)
*Fall Semester 2024*
**Teaching Assistant** at the School of Data Science, Fudan University, with Prof. Jiaqing Liang

- Responsible for all lecture materials related to Operating Systems: Process, Virtural Memory, and System-level I/O. Prepared relevant lecture slides and lab tutorials, and Shell Lab project.

- With the ever-increasing importance of system infrastructure knowledge in LLM research, aimed to lay the foundation of system education for future NLP researchers, and guide them into the fascinating world of LLMs and MLSys.

## HONORS & AWARDS

**Deans' Honors List (Top 2% in the College of Letters and Science)**
*University of California, Davis, Fall Quarter 2023*

**Academic Excellence Scholarship (Ranked first in the major)**
*Fudan University, 2022 - 2023*

**Panasonic Scholarship (Ranked first in the major)**
*Fudan University, 2021 - 2022*

## SKILLS

**Programming:** C/C++, Python, MATLAB, R, HTML/CSS/JavaScript, SQL, LaTeX
**Frameworks:** Pytorch, HuggingFace, OpenMP/MPI
**Languages:** English (TOEFL iBT 107, speaking 24), Chinese (native)