

JIAJUN ZHU

3200106048@zju.edu.cn · (+86) 173-009-89120 · jurray-jiajun.github.io

EDUCATION

Zhejiang University

Sept. 2020 - Now

B.S. in Mathematics

- Major GPA: 3.80/4.00 (top 1/64 in sophomore year)
- Math Average: 91.3
- Core Courses:
 - Ordinary Differential Equation: 95
 - Mathematical Software: 99
 - Advanced Algebra II: 95
 - Mathematical Analysis II/III: 93
 - Mathematical Statistics: 92
 - Abstract Algebra: 90
 - Real Variable Analysis: 93
 - Point Set Topology: 93
 - Geometry: 90

PUBLICATION

1. **Jiajun Zhu**, Peihao Wang, Ruisi Cai, Jason D. Lee, Pan Li, Zhangyang Wang. Rethinking Addressing in Language Models via Contextualized Equivariant Positional Encoding, *The Thirteenth International Conference on Learning Representations (ICLR)*, 2025. Submitted.
2. **Jiajun Zhu**, Siqi Miao, Rex Ying, Pan Li. Towards Understanding Sensitive and Decisive Patterns in Explainable AI: A Case Study of Model Interpretation in Geometric Deep Learning, *Nature Machine Intelligence*, 2024. Minor revision. (Manuscript is deposited on arXiv)
3. Peihao Wang, Ruisi Cai, Yuehao Wang, **Jiajun Zhu**, Pragya Srivastava, Zhangyang Wang, Pan Li, Understanding Bottlenecks of State Space Models through the Lens of Recency and Over-smoothing, *The Thirteenth International Conference on Learning Representations (ICLR)*, 2025. Submitted.
4. Yifei Sun, Qi Zhu, Yang Yang, Chunping Wang, Tianyu Fan, **Jiajun Zhu**, Lei Chen. Fine-tuning Graph Neural Networks by Preserving Graph Generative Patterns, *The Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI)*, 2024.

SELECTED RESEARCH

Enhanced Transformer with Contextualized Equivariant Position Encoding

May. 2024 - Sept. 2024

Supervised by Prof. Zhangyang Wang

University of Texas at Austin

- Introduced TAPE, a novel framework for dynamically layer-updated positional embeddings in transformers, adapting to content and surpassing the limitations of fixed long-term decay in traditional positional embeddings.
- Proposed the principles of permutation invariance and orthogonal equivariance to enhance the generalization of positional embeddings, and designed an enhanced Transformer with modules that integrate positional information into both the attention and feedforward layers in alignment with these principles.
- Demonstrated that TAPE excels in language modeling and downstream tasks such as arithmetic reasoning and long-context retrieval, achieving strong performance in both pretraining from scratch and parameter-efficient fine-tuning.

Interpretability of Geometric Deep Learning for Scientific Tasks

Dec. 2022 - Apr. 2024

Supervised by Prof. Pan Li

Georgia Institute of Technology

- Adapted 12 interpretability techniques from graph neural networks to geometric deep learning models, which are widely employed in scientific tasks, and benchmarked their performance.
- Proposed the definition of two critical concepts in the domain of interpretability: sensitive patterns and decisive patterns, highlighting their misalignment, an aspect previously overlooked by researchers.
- Derived key insights from empirical evidence to guide the effective and appropriate application of two major categories of interpretability techniques: post-hoc methods and self-interpretable methods.

ONGOING RESEARCH

Interpretable State-Space Model with Enhanced Selective Mechanism

Oct. 2024 - Now

Supervised by Prof. Zhangyang Wang

University of Texas at Austin

- Proposed multiple approaches to implement an interpretable version of Mamba, named S7, leveraging the information bottleneck technique.
- Demonstrated that S7 achieves inherent interpretability without compromising generation performance, as validated on common language benchmarks.

Beyond Linear Separability of Attack-Defense in Latent Space of LLMs

Dec. 2024 - Now

Supervised by Prof. Prateek Mittal

Princeton University

- Proposed a backdoor attack that breaks the linear separability of attack-defense by injecting diversified poisoned data, enabling it to bypass defenses like BEEAR and Refusal-Direction based on linear separability assumptions.

FULL EXPERIENCE

LLM-based Agent for Automatic Cell Type Annotations

May. 2024 - Dec. 2024

Remote Research Intern Supervised by Prof. Zhiting Hu

University of California San Diego

- Proposed a framework empowering LLM-based agents to generate hypotheses, conduct experiments, perform evaluations, and iteratively refine hypotheses based on evaluation outcomes.
- Validated the approach through benchmarking on cell type annotation datasets and conducted a user study utilizing a streamlined user interface designed for scientists.

Interactive Reasoning of Visual Language Models

Dec. 2023 - Apr. 2024

Research Intern Supervised by Prof. Yaochu Jin

Westlake University

- Proposed a paradigm enabling interaction with visual language models (VLMs) through visual referencing, specifically utilizing “click and segment” actions to improve interactivity and reference accuracy.
- Enhanced the reasoning capabilities of VLMs for image segmentation by fine-tuning them on a custom-built multi-modal dataset.

Mitigating Structural Divergence in Fine-tuning Graph Neural Networks

Aug. 2022 - Oct. 2022

Supervised by Prof. Yang Yang

Zhejiang University

- Provided a theoretical analysis using Taylor decomposition to guide the design of our method, breaking down the module output from a complete matrix into a linear combination of coefficients and bases.
- Validated the effectiveness of our method by implementing five baseline approaches and benchmarking performance across eight molecular datasets.

Graph Neural Networks for Electronic Property Prediction

Jan. 2022 - Jun. 2022

Team Leader Supervised by Prof. Renjun Xu

Zhejiang University

- Collected 100k+ data of electronic density and band structure from material database.
- Adapted CGCNN and MEGNet models for accurate electronic property prediction.

ADDITIONAL INFORMATION

- Programming languages: Python, C/C++, CUDA.
- Software & Frameworks: LaTeX, Git, PyTorch, PyTorch Geometric, Transformers, Triton, Equinox.
- Interests: Guitar (performed at the school’s New Year party), Skiing, Basketball, Snooker.