

Runlong Zhou (周润龙)

Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, WA 98195, USA  
zhourunlongvector@gmail.com, vectorzh@cs.washington.edu  
vectorzhou.com

February 17, 2025

Research Interests

- Reinforcement learning (RL) + large language models (LLMs)
- RL theory

Education

- **University of Washington**  
*PhD, Paul G. Allen School of Computer Science & Engineering, advised by Prof. Simon S. Du* Seattle, USA  
2022.9 - Now (Est. 2027.8)
- **Tsinghua University**  
*BEng, Special Pilot CS Class (Yaoclass), Institute for Interdisciplinary Information Sciences (IIIS)* Beijing, China  
*GPA: 3.84 (overall) / 3.89 (major) over scale 4.0, Rank: 16 / 54* 2018.8 - 2022.6

Academic Experience

- **Microsoft Research**  
*Research intern with Dr. Baolin Peng* Redmond, USA  
2025.2 - Now
- **Microsoft AI**  
*Research intern with Dr. Yi Zhang* Redmond, USA  
2024.6 - 2024.8
- **Microsoft Research**  
*Research intern with Dr. Beibin Li* Redmond, USA  
2023.6 - 2023.9
- **University of Washington**  
*Research intern with Prof. Simon S. Du* Virtual  
2020.9 - 2022.9
- **Facebook AI Research**  
*Research intern with Dr. Alessandro Lazaric and Dr. Matteo Pirodda* Virtual  
2021.3 - 2021.5

Publications

\* denotes equal contribution or alphabetical ordering.

1. The Crucial Role of Samplers in Online Direct Preference Optimization [Link]  
Ruizhe Shi\*, Runlong Zhou\*, Simon S. Du  
ICLR 2025  
We prove that online DPO with a mixture of samplers achieves quadratic convergence with exact gradients and linear convergence with estimations.

Poster

2. Reflect-RL: Two-Player Online RL Fine-Tuning for LMs [Link]  
Runlong Zhou, Simon S. Du, Beibin Li  
ACL 2024  
We developed Reflect-RL, a two-player system to align language models with interactive decision-making tasks. Techniques included are reflection, negative example generation, single-prompt action enumeration, and curriculum learning.

Poster

3. **Free from Bellman Completeness: Trajectory Stitching via Model-based Return-conditioned Supervised Learning** [Link]  
Zhaoyi Zhou, Chuning Zhu, **Runlong Zhou**, Qiwen Cui, Abhishek Gupta, Simon S. Du  
ICLR 2024 Poster  
*We study the strengths and weaknesses of return-conditioned supervised learning, and propose an empirically improved algorithm.*
4. **Sharp Variance-Dependent Bounds in Reinforcement Learning: Best of Both Worlds in Stochastic and Deterministic Environments** [Link]  
**Runlong Zhou**, Zihan Zhang, Simon S. Du  
ICML 2023 Poster  
*We provide a systematic study of variance-dependent regret bounds of model-based and model-free reinforcement learning for tabular MDPs. The proposed model-based algorithm is both optimal for stochastic and deterministic MDPs.*
5. **Variance-Dependent and Horizon-Free Reinforcement Learning for Latent Markov Decision Processes** [Link]  
**Runlong Zhou**, Ruosong Wang, Simon S. Du  
ICML 2023 Poster  
*We provide an algorithm framework for Latent MDPs (with context in hindsight), achieving the first horizon-free minimax regret. We complement the study by giving a novel regret lower bound for LMDPs using the symmetrization technique.*
6. **Understanding Curriculum Learning in Policy Optimization for Online Combinatorial Optimization** [Link]  
**Runlong Zhou**, Zelin He, Yuandong Tian, Yi Wu, Simon S. Du  
TMLR  
*We formulate of canonical online Combinatorial Optimization problems as Latent MDPs and give convergence guarantee of Natural Policy Gradient on LMDPs. We show effectiveness of Curriculum Learning through the perspective of relative conditional number.*
7. **Stochastic Shortest Path: Minimax, Parameter-Free and Towards Horizon-Free Regret** [Link]  
Jean Tarbouriech\*, **Runlong Zhou\***, Simon S. Du, Matteo Pirodda, Michal Valko, Alessandro Lazaric  
NeurIPS 2021 Spotlight, 3% acceptance rate  
*We propose an algorithm (EB-SSP) for SSP problems, which is the first to achieve minimax optimal regret while being parameter-free.*

## Preprints

\* denotes equal contribution or alphabetical ordering.

1. **Transformers are Efficient Compilers, Provably** [Link]  
Xiyu Zhai, **Runlong Zhou**, Liao Zhang, Simon S. Du  
*We propose Cybertron as a proof vehicle for transformers' expressive ability and show that for a compilation task, transformers need only a logarithm number of parameters while any recurrent neural network needs at least a linear number of parameters.*
2. **Multi-Agent Reinforcement Learning from Human Feedback: Data Coverage and Algorithmic Techniques** [Link]  
Natalia Zhang\*, Xinqi Wang\*, Qiwen Cui\*, **Runlong Zhou**, Sham M. Kakade, Simon S. Du  
*We study Multi-Agent Reinforcement Learning from Human Feedback (MARLHF) by exploring both theoretical foundations and empirical validations. Included in our proposed methods are Mean Squared Error (MSE) regularization and imitation learning.*

## Ongoing Projects

1. **Nash Learning from Human Feedback**  
Advised by Prof. Simon S. Du  
*We study the Nash equilibrium of the two-player constant-sum game motivated by reinforcement learning from human feedback (RLHF) instances. We aim to give algorithms with better provable convergence rates and adapt them to the training of large language models (LLMs) for better efficiency.*
2. **Physics of Language Models**  
Advised by Dr. Yi Zhang  
*We study how language models learn and manipulate knowledge.*

## Awards, Grants & Honors

- Undergraduate:

IIIS Outstanding Graduate . . . . .	2022
The 2021 China Collegiate Programming Contest, Guilin Site (Gold Medal) . . . . .	2021
IIIS Research Innovation Scholarship . . . . .	2021
IIIS Academic Performance Scholarship . . . . .	2021
Tsinghua University Air Rifle Competition (First Place) . . . . .	2021
The 2019 ACM-ICPC Asia Regional Contest, Xuzhou Site (Gold Medal) . . . . .	2019
The 2018 ACM-ICPC Asia Regional Contest, Beijing Site (Gold Medal) . . . . .	2018

- Secondary school:

The 34th National Olympiad in Informatics (Silver Medal) . . . . .	2017
China Team Selection Competition 2017 (Gold Medal) . . . . .	2017
The 2016 ACM-ICPC Asia CHINA-Final Contest (Gold Medal) . . . . .	2016
The 2016 China Collegiate Programming Contest Finals (Silver Medal) . . . . .	2016
The 33rd National Olympiad in Informatics (Silver Medal) . . . . .	2016

## Other Projects

1. Ray Tracing Renderer [Link]

Runlong Zhou

*Advanced Computer Graphics course project*

*Optimized path tracing framework supporting Mitsuba configurations, many textures and sampling methods*

C++

## Miscellanea

- **Professional skills:** Algorithm design, Data structures, Deep (Reinforcement) learning
- **Programming skills:** C++ / C, python,  $\text{\LaTeX}$ , HTML, JavaScript (Node.js), CUDA, Java, MATLAB, Rust
- **Hobbies:** Air rifle / pistol shooting, Archery
- **Teaching:** Teaching Olympiad in Informatics to secondary school students between 2018 and 2021. Teaching assistant of UW CSE 543 Deep Learning, CSE 446 Machine Learning.
- **Reviewing:** ICML 2022, NeurIPS 2022, ICML 2024, COLT 2024, ICLR 2025.