

Education

KAIST (Korea Advanced Institute of Science and Technology)

MASTER'S DEGREE IN GRADUATE SCHOOL OF ARTIFICIAL INTELLIGENCE

- GPA: 4.2 / 4.3 (~ 3rd semester), Supervised by Professor James Thorne

Seoul, S.Korea

Sep. 2023 - Aug. 2025 (Expected)

SKKU (SungKyunKwan University)

BACHELOR'S DEGREE IN STATISTICS & INDUSTRIAL ENGINEERING

- GPA : 4.38 / 4.5 (Summa Cum Laude)

Seoul, S.Korea

Mar. 2017 - Feb. 2023

Research Interest

Natural Language Processing

Large Language Model, Preference Alignment, Reasoning

Reinforcement Learning

Reinforcement Learning with Human/AI Feedback/Verifiable Rewards

Publications

2025

On the Robustness of Reward Models for Language Model Alignment

Jiwoo Hong, Noah Lee, EunKi Kim, Guijin Son, Woojin Chung, Aman Gupta, Shao Tang, & James Thorne

ICML 2025

Released: 25.05.12, Citation: -

AlphaPO - Reward Shape Matters for LLM Alignment

Aman Gupta, Shao Tang, Qingquan Song, Siyou Zhu, Jiwoo Hong, Ankan Saha, Viral Gupta, Noah Lee, EunKi Kim, Siyu Zhu, Parag Agrawal, Natesh Pillai, & S. Sathya Keerthi

ICML 2025

Released: 24.01.07, Citation: 2

Linguistic Generalizability of Test-Time Scaling in Mathematical Reasoning

Guijin Son, Jiwoo Hong, Hyunwoo Ko, & James Thorne

ACL 2025

Released: 25.02.25, Citation: 4

Online Difficulty Filtering for Reasoning Oriented Reinforcement Learning

Sanghwan Bae*, Jiwoo Hong*, Min Young Lee, Hanbyul Kim, Jeongyeon Nam, & Donghyun Kwak

Under Review for COLM 2025

Released: 25.04.04, Citation: 2

Cross-lingual Transfer of Reward Models in Multilingual Alignment

Jiwoo Hong*, Noah Lee*, Rodrigo Martínez-Castaño, César Rodríguez, & James Thorne

NAACL 2025

Released: 24.10.24, Citation: 3

2024

ORPO: Monolithic Preference Optimization without Reference Model

Jiwoo Hong, Noah Lee & James Thorne

EMNLP 2024

Released: 24.03.12, Citation: 302

Stable Language Model Pre-training by Reducing Embedding Variability

Woojin Chung, Jiwoo Hong, Na Min An, James Thorne, & Se Young Yoon

EMNLP 2024

Released: 24.09.12, Citation: 2

Margin-aware Preference Optimization for Aligning Diffusion Models without Reference

Jiwoo Hong*, Sayak Paul*, Noah Lee, Kashif Rasul, James Thorne & Jongheon Jeong

ICLR 2025 SCOPE Workshop

Released: 24.06.11, Citation: 7

Evaluating the Consistency of LLM Evaluators

Noah Lee*, Jiwoo Hong*, & James Thorne

COLING 2025

Released: 24.11.30, Citation: 6

2023

Disentangling Structure and Style: Political Bias Detection in News by Inducing Document Hierarchy

Jiwoo Hong, Yejin Cho, Jaemin Jung, Jiyoung Han & James Thorne

Findings of EMNLP 2023

Released: 23.12.06, **Citation:** 7

2022

MARL-Based Dual Reward Model on Segmented Actions for Multiple Mobile Robots in Automated Warehouse Environment

Hyeoksoo Lee, Jiwoo Hong & Jongpil Jeong

Applied Science

Released: 22.05.07, **Citation:** 11

Industry Experiences

Amazon - Rufus

APPLIED SCIENTIST INTERNSHIP

Palo Alto, CA

May. 2025 - Present

- **Topic:** Multi-objective optimization in language model post-training with reinforcement learning.

Naver Cloud

RESEARCH INTERNSHIP

Seoul, S.Korea

Feb. 2025 - May. 2025

- **Topic:** Reinforcement learning with verifiable rewards (RLVR) algorithms and pipeline development.

Additional Experiences

Multiple Invited Talks on LLM Alignment

Seoul, S.Korea

INVITED SPEAKER

- Talk on recent alignment techniques and ORPO at:
- (1) Kakao Brain, (2) AI EXPO Korea, (3) KISTI (Korean Institute of Science and Technology Information), (4) Twelve Labs,
- and (5) Trillion Parameter Consortium Seminar.

Zephyr-ORPO Project

Seoul, S.Korea

OPEN-SOURCE LANGUAGE MODEL DEVELOPMENT

April. 2024

- Developing data-efficient open-source instruction-following large language models with fine-grained synthetic data.
- In collaboration with Argilla and Hugging Face.

ML/NLP Conference Reviewer

Seoul, S.Korea

REVIEWER

Nov. 2024 - Present

- Served as a reviewer for ICLR 2025 on the papers related to LLM alignment.
- Served as a reviewer for ACL Rolling Review (ARR) February cycle on the papers.