

Jiwoo Hong

✉ jiwoo_hong@kaist.ac.kr | 📄 jiwooya1000 | 📄 jiwoohong09 | 🐦 @jiwoohong98 | 📄 Google Scholar | 🌐 Website

Mountain View, CA, USA

Industry Experiences

LinkedIn

SPECIALIST AI ENGINEER, CORE AI

- Agentic AI system development and evaluation for LinkedIn.

Mountain View, CA, US

Jan. 2026 - Present

Amazon

APPLIED SCIENTIST INTERN

- Steerability in LLMs via multi-objective reinforcement learning (MORL) for Amazon Rufus.

Palo Alto, CA, US

May. 2025 - Aug. 2025

Naver Cloud

AI RESEARCH SCIENTIST INTERN

- Core contributor in developing a proprietary large-scale reasoning model, HyperCLOVA X THINK.

Seoul, S.Korea

Feb. 2025 - May. 2025

Education

KAIST (Korea Advanced Institute of Science and Technology)

MASTER'S DEGREE IN GRADUATE SCHOOL OF ARTIFICIAL INTELLIGENCE

- GPA: 4.12 / 4.3 (Supervised by Professor James Throne)

Seoul, S.Korea

Sep. 2023 - Aug. 2025

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

Publications and Patents

Selected Works

ORPO: Monolithic Preference Optimization without Reference Model

Jiwoo Hong, Noah Lee & James Thorne

EMNLP 2024

Topic: Offline Preference Learning

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

Topic: RLHF, Reward Modeling

Online Difficulty Filtering for Reasoning Oriented Reinforcement Learning

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

EACL 2026

Topic: RLVR, Reasoning

2026

Bayesian Preference Learning for Test-Time Steerable Reward Models

Jiwoo Hong*, Shao Tang*, & Zhipeng Wang

Under Review for ICML 2026

Topic: RLHR, Bayesian Learning

Online Difficulty Filtering for Reasoning Oriented Reinforcement Learning

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

EACL 2026

Topic: RLVR, Reasoning

Margin-aware Preference Optimization for Aligning Diffusion Models without Reference

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

AAAI 2026

Topic: RLHF, Diffusion

2025

Method and Apparatus for Voice Profiling

Jin Yeong Bak*, Jiwoo Hong*, Seung Woo Lee*, Hyolim Jeon*, Soo Jeong Lee* & Se Hyun Ahn*

KR Patent (ID: 102883231)

Topic: Deep Learning

HyperCLOVA X THINK Technical Report
NAVER CLOUD HYPERCLOVA X TEAM (CORE CONTRIBUTOR)

Technical Report
Topic: Proprietary Reasoning Model

When AI Co-Scientists Fail: SPOT-a Benchmark for Automated Verification of Scientific Research

Preprint

GUIJIN SON, **Jiwoo Hong**, HONGLU FAN, HEEJEONG NAM, HYUNWOO KO, SEUNGWON LIM, JINYEOP SONG, JINHA CHOI, GONÇALO PAULO, YOUNGJAE YU, & STELLA BIDERMAN

Topic: AI for Science

On the Robustness of Reward Models for Language Model Alignment

ICML 2025

Jiwoo Hong, NOAH LEE, EUNKI KIM, GUIJIN SON, WOOJIN CHUNG, AMAN GUPTA, SHAO TANG, & JAMES THORNE

Topic: RLHF, Generalizability

AlphaPO: Reward Shape Matters for LLM Alignment

ICML 2025

AMAN GUPTA, SHAO TANG, QINGQUAN SONG, SIROU ZHU, **Jiwoo Hong**, ANKAN SAHA, VIRAL GUPTA, NOAH LEE, EUNKI KIM, SIYU ZHU, PARAG AGRAWAL, NATESH PILLAI, & S. SATHIYA KEERTHI

Topic: RLHF, Interpretability

Linguistic Generalizability of Test-Time Scaling in Mathematical Reasoning

ACL 2025

GUIJIN SON, **Jiwoo Hong**, HYUNWOO KO, & JAMES THORNE

Topic: Reasoning, Multilingual

Cross-lingual Transfer of Reward Models in Multilingual Alignment

NAACL 2025

Jiwoo Hong*, NOAH LEE*, RODRIGO MARTÍNEZ-CASTAÑO, CÉSAR RODRÍGUEZ, & JAMES THORNE

Topic: RLHF, Generalizability

2024

ORPO: Monolithic Preference Optimization without Reference Model

EMNLP 2024

Jiwoo Hong, NOAH LEE & JAMES THORNE

Topic: Offline Preference Learning

Stable Language Model Pre-training by Reducing Embedding Variability

EMNLP 2024

WOOJIN CHUNG, **Jiwoo Hong**, NA MIN AN, JAMES THORNE, & SE YOUNG YOON

Topic: Pre-training

Evaluating the Consistency of LLM Evaluators

COLING 2025

NOAH LEE*, **Jiwoo Hong***, & JAMES THORNE

Topic: LLM-as-a-Judge

2023

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

Findings of EMNLP 2023

Jiwoo Hong, YEJIN CHO, JAEMIN JUNG, JIYOUNG HAN & JAMES THORNE

Topic: NLP Application

2022

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

Applied Science

HYEOKSOO LEE, **Jiwoo Hong** & JONGPIL JEONG

Topic: Multi-agent RL

Additional Experiences

Multiple Invited Talks on LLM Alignment and Post-training

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) Tweleve Labs,
- and (5) Trillion Parameter Consortium Seminar.

Zephyr-ORPO Project

Seoul, S.Korea

OPEN-SOURCE LANGUAGE MODEL DEVELOPMENT

April. 2024

- Developing data-efficient open-recipe instruction-following LLMs with fine-grained synthetic data and Mixtral-8x22B-v0.1.
- In collaboration with Argilla and Hugging Face, published Zephyr-ORPO-141B-A35B-v0.1.

HyperCLOVA X THINK

Seoul, S.Korea

RLVR PIPELINE DEVELOPMENT IN PROPRIETARY LARGE REASONING MODEL

Feb.-May. 2025

- Participated in developing a proprietary large-scale reasoning model, HyperCLOVA X THINK, specialized in Korean.
- Presented theoretical background on online difficulty filtering in reinforcement learning with verifiable rewards (RLVR), contributing to sample-efficient improvements in complex reasoning tasks.

ML/NLP Conference Reviewer

-

REVIEWER

Nov. 2024 - Present

- ICML 2026
- AAAI 2026
- ICLR 2025-2026
- NeurIPS 2026
- TMLR 2025-2026
- ACL Rolling Review 2025-2026