





JISEON KIM

Ph.D. candidate

 jiseon\_kim@kaist.ac.kr  hikoseon12.github.io  
 jiseon-kim-8ab574136  github.com/hikoseon12

SUMMARY

KAIST Ph.D. candidate specializing in AI Safety and AI for Social Good ●**AI Safety**: Evaluating cultural bias and moral reasoning in LLMs ●**AI for Social Good**: Developing scalable AI frameworks for domain-specific adaptation, policy analytics ●**Skills**: LLM evaluation and benchmarking, dataset construction, domain adaptation of AI models, and data-driven insights

EDUCATION

3/2020-Present	<b>Korea Advanced Institute of Science and Technology (KAIST)</b> Ph.D. candidate in School of Computing	Daejeon, Korea
3/2017-2/2019	<b>Korea Advanced Institute of Science and Technology (KAIST)</b> Master in School of Computing	Daejeon, Korea
3/2013-2/2017	<b>Sookmyung Women's University</b> B.S. in Computer Science Graduated with the highest honor (1/68)	Seoul, Korea

EXPERIENCE

8/2025 - 10/2025	<b>Research Intern @ Max Planck Institute for Security and Privacy (MPI-SP)</b> • Conducted research on evaluating LLM alignment by analyzing how models understand diverse human personas in moral dilemmas	MPI-SP
5/2024 - 5/2024 7/2022 - 8/2022 6/2019 - 8/2019	<b>Visiting Researcher @ Massachusetts Institute of Technology (MIT)</b> • Conducted interdisciplinary research with political science and developed a scalable AI framework to understand the U.S. legislative process • Work published at EMNLP 2021 - "Learning Bill Similarity with Annotated and Augmented Corpora of Bills" • Preprinted work (2025) on "Measuring Interest Group Positions on Legislation: An AI-Driven Analysis of Lobbying Reports."	MIT
3/2023 - 6/2023	<b>Research Intern @ NAVER AI Lab</b> • Constructed a Korean bias benchmark dataset to make safer and trustworthy Korean LLM • Work published at TACL 2024 - "KoBBQ: Korean Bias Benchmark for Question Answering"	NAVER AI Lab

PUBLICATION

Preprint 2025	<b>[P9] Measuring Interest Group Positions on Legislation: An AI-Driven Analysis of Lobbying Reports</b> MIT Jiseon Kim, Dongkwan Kim, In Song Kim, Alice Oh • Expanded lobbying position classification beyond binary to include nuanced categories • Built scalable AI framework with LLMs and GNNs to annotate 279K+ interest group-bill pairs and compute policy preference scores • Analyzed lobbying strategies influenced by policy area, legislative stage, and group size	
Preprint 2024	<b>[P8] Uncovering Factor Level Preferences to Improve Human-Model Alignment</b> Juhyun Oh*, Eunsu Kim*, Jiseon Kim, Wenda Xu, Inha Cha, William Yang Wang, Alice Oh <sup>†</sup> (equal contribution) • Developed PROFILE, a framework to explain factors driving LLM-human preference alignment • Identified key differences in preferences between humans and LLMs across tasks • Emphasized explainable analysis to enhance human-model alignment and training	
EMNLP 2024 Long paper	<b>[P7] Perceptions to Beliefs: Exploring Precursory Inferences for Theory of Mind in Large Language Models</b> Chani Jung, Dongkwan Kim, Jiho Jin, Jiseon Kim, Yeon Seonwoo, Yejin Choi, Alice Oh, Hyunwoo Kim • Introduced Percept-ToMi and Percept-FANToM datasets to assess ToM precursors in LLMs • Demonstrated LLMs excel in perception inference but show limitations in perception-to-belief inference • Developed PercepToM, a method that improves LLM performance on ToM benchmarks	Allen AI
Technical Report 2024	<b>[P6] HyperCLOVA X Technical Report</b> Kang Min Yoo et al., Jiseon Kim, ... • Introduced LLM optimized for Korean language and culture, with strong English, math, and coding skills • Trained on Korean, English, and code data, and evaluated on various benchmarks in both languages • Contributed to model evaluations, including bias measurement in Korean culture through KoBBQ	NAVER AI Lab

TACL 2024, present at ACL 2024	<b>[P5] KoBBQ: Korean Bias Benchmark for Question Answering</b> <span>NAVER AI Lab</span> Jiho Jin*, Jiseon Kim*, Nayeon Lee*, Hanual Yoo*, Alice Oh, Hwaran Lee <sup>†</sup> (equal contribution) <ul style="list-style-type: none"> <li>Introduced a Korean bias benchmark dataset to address challenges in adapting to non-US cultures</li> <li>Proposed a framework for cultural adaptation, categorizing and validating biases via a large-scale survey</li> <li>Revealed significant differences in LM biases compared to a machine-translated version, highlighting the need for culturally-sensitive benchmarks</li> </ul>
EMNLP 2021 Long paper	<b>[P4] Learning Bill Similarity with Annotated and Augmented Corpora of Bills</b> <span>MIT</span> Jiseon Kim, Elden Griggs, In Song Kim, Alice Oh <ul style="list-style-type: none"> <li>Proposed a 5-class task for bill document semantic similarities to understand bill-to-bill linkage in the legislative process</li> <li>Improved model performance by achieving a 5.5% higher F1 score compared to the baseline using data augmentation and multi-stage training</li> <li>Quantified the similarities across legal documents at various levels of aggregation</li> </ul>
EMNLP 2021 Short paper	<b>[P3] Efficient Contrastive Learning via Novel Data Augmentation and Curriculum Learning</b> Seonghyeon Ye, Jiseon Kim, Alice Oh <ul style="list-style-type: none"> <li>Proposed a memory-efficient continual pretraining method</li> <li>Outperformed baseline models on GLUE benchmark with only 70% computational memory usage</li> </ul>
EMNLP 2021 Long paper	<b>[P2] Dimensional emotion detection from categorical emotion</b> Sungjoon Park, Jiseon Kim, Seonghyeon Ye, Jaeyeol Jeon, Hee Young Park, Alice Oh <ul style="list-style-type: none"> <li>Utilized categorical emotion annotations to train a model predicting fine-grained emotions</li> <li>Optimized model with Earth Mover's Distance loss to predict fine-grained and categorical emotions</li> <li>Achieved comparable performance to state-of-the-art classifiers in emotion classification</li> </ul>
IEEE transactions on intelligent transportation systems 2020	<b>[P1] Denoising recurrent neural networks for classifying crash-related events</b> Sungjoon Park, Yeon Seonwoo, Jiseon Kim, Jooyeon Kim, Alice Oh <ul style="list-style-type: none"> <li>Developed efficient neural network model with noisy time-series data with missing values for crash event classification</li> <li>Outperformed baseline models, improving event classification accuracy in driving scenarios</li> </ul>

## WORKSHOP

BiAlign @ICLR 2025	<b>[W3] Exploring Persona-dependent LLM Alignment for the Moral Machine Experiment</b> <span>MPI-SP</span> Jiseon Kim*, Jea Kwon*, Luiz Felipe Vecchietti*, Alice Oh, Meeyoung Cha <sup>†</sup> (equal contribution)
WiML @NeurIPS 2024	<b>[W2] Understanding Lobbying Strategies in Legislative Process: Bill Position Dataset and Lobbying Analysis</b> <span>MIT</span> Jiseon Kim, Dongkwan Kim, Joohye Jeong, In Song Kim, Alice Oh
C3NLP @ACL 2024	<b>[W1] KoBBQ: Korean Bias Benchmark for Question Answering</b> <span>NAVER AI Lab</span> Jiho Jin*, Jiseon Kim*, Nayeon Lee*, Hanual Yoo*, Alice Oh, Hwaran Lee <sup>†</sup> (equal contribution)

## INVITED TALK

ExploreCSR@Google March 21, 2025	<b>Uncovering the Hidden Politics of Lawmaking: How Bills and Lobbying Shape U.S. Policy</b> <span>KAIST</span> Presented on AI for Political Science to understand the legislative process, supported by Google and hosted by KAIST School of Computing.
MPI-SP@Germany Feb 25, 2025	<b>LLMs and the Political-Cultural Lens in Social Science</b> <span>MPI-SP</span> Invited talk at Max Planck Institute for Security and Privacy, hosted by Prof. Meeyoung Cha (Data Science for Humanity).
MLAI@Yonsei Jan 2, 2025	<b>Things I Wish I Had Known Earlier in Grad School</b> <span>Yonsei University</span> Invited talk on networking, self-promotion, and collaboration in academia, hosted by Prof. Kyungwoo Song at the Machine Learning and Artificial Intelligence (MLAI) Lab.

## AWARD, SCHOLARSHIP & FUNDING

4/2025 - Present 12/2019 - 8/2024	<b>MISTI Global Seed Funds</b> <span>MIT</span> MIT's Global Seed Funds facilitate international collaborations for addressing global challenges
10/2024	<b>2024 KAIST Graduate Student Outstanding Paper Award</b> <span>KAIST</span> Awarded for KoBBQ: Korean Bias Benchmark for Question Answering
3/2020 - Present	<b>KAIST Support Scholarship (Ph.D.)</b> <span>KAIST</span>