

Email: [dg.lee@postech.ac.kr](mailto:dg.lee@postech.ac.kr)    Homepage: <https://donggeon.github.io>    Google Scholar: [/DongGeon Lee](#)

RESEARCH INTERESTS	Data-centric natural language processing (NLP). Trustworthy and safe Large Language Models (LLMs): safety evaluations, guardrails, and automated red teaming.	
EDUCATION	<b>M.S. student in Artificial Intelligence</b> <i>Pohang University of Science and Technology (POSTECH)</i>	Feb 2024 - Present <i>Pohang, South Korea</i>
	<b>B.S. in Information and Communication Engineering</b> <i>Inha University</i>	Mar 2018 - Feb 2024 <i>Incheon, South Korea</i>
RESEARCH EXPERIENCES	<b>Graduate Research Assistant</b> <i>Data Intelligence Lab, POSTECH (Advisor: Prof. Hwanjo Yu)</i>	Feb 2024 - Present <i>Pohang, South Korea</i>
	<ul style="list-style-type: none"><li>• Research on Vision-Language Model safety benchmarks and evaluation methodologies.</li><li>• Research on knowledge conflicts of LLMs between external and internal knowledge.</li></ul>	
	<b>Research Scientist</b> (Freelance) <i>AIM Intelligence</i>	Jul 2025 - Present <i>Seoul, South Korea</i>
	<ul style="list-style-type: none"><li>• Research on automated multi-turn red teaming for LLMs.</li><li>• Research on safety guardrails for multimodal/multilingual LLMs.</li><li>• Research on dynamic evaluation of LLM compliance with organization-defined policies.</li></ul>	
	<b>Research Intern</b> <i>KT Corporation</i>	Jan 2025 - Feb 2025 <i>Seoul, South Korea</i>
INTERNATIONAL PUBLICATIONS	<b>Selected Publications</b> (Full list available on <a href="#">Google Scholar</a> .) * Equal contribution; † Equal advising	
	<ul style="list-style-type: none"><li>[6] COMPASS: A Framework for Evaluating Organization-Specific Policy Alignment in LLMs Dasol Choi*, <u>DongGeon Lee</u>*, Brigitta Jesica Kartono*, Helena Berndt, Haon Park, Hwanjo Yu†, Minsuk Kahng† Under Review, 2025.10</li><li>[5] Are Vision-Language Models Safe in the Wild? A Meme-Based Benchmark Study <u>DongGeon Lee</u>*, Joonwon Jang*, Jihae Jeong, Hwanjo Yu <a href="#">EMNLP'25</a>   The 2025 Conference on Empirical Methods in Natural Language Processing</li><li>[4] Everyday Physics in Korean Contexts: A Culturally Grounded Physical Reasoning Benchmark Jihae Jeong*, DaeYeop Lee*, <u>DongGeon Lee</u>, Hwanjo Yu <a href="#">MRL @ EMNLP'25</a>   5th Multilingual Representation Learning Workshop (Co-located with the 2025 Conference on Empirical Methods in Natural Language Processing)</li><li>[3] When Good Sounds Go Adversarial: Jailbreaking Audio-Language Models with Benign Inputs Bodam Kim*, Hiskias Dingeto*, Taeyoun Kwon*, Dasol Choi, <u>DongGeon Lee</u>, Haon Park, Jae-Hoon Lee, Jongho Shin <a href="#">arXiv Preprint</a>, 2025.08</li></ul>	

- [2] Typed-RAG: Type-Aware Decomposition of Non-Factoid Questions for Retrieval-Augmented Generation  
 DongGeon Lee\*, Ahjeong Park\*, Hyeri Lee, Hyeonseo Nam, Yunho Maeng  
[XLLM @ ACL'25](#) | The First Workshop on Structure-aware Large Language Models (Co-located with the 63rd Annual Meeting of the Association for Computational Linguistics)  
[NAACL'25 SRW](#) (Non-Archival) | 2025 Annual Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Student Research Workshop
- [1] REFIND at SemEval-2025 Task 3: Retrieval-Augmented Factuality Hallucination Detection in Large Language Models  
 DongGeon Lee, Hwanjo Yu  
[SemEval @ ACL'25](#) | The 19th International Workshop on Semantic Evaluation (Co-located with the 63rd Annual Meeting of the Association for Computational Linguistics)

## ACADEMIC SERVICES

Reviewer of AAAI'25 (The Association for the Advancement of Artificial Intelligence)	2025
Reviewer of International Journal of Nursing Studies	2025
Reviewer of MELT (Workshop on Multilingual and Equitable Language Technologies) at COLM'25	2025
Reviewer of SemEval (International Workshop on Semantic Evaluation) at ACL'25	2025
Student Volunteer of ACL'25 (Annual Meeting of the Association for Computational Linguistics)	2025

## HONORS AND AWARDS

Excellent Paper Award <i>HCLT 2025 (The 37th Annual Conference on Human &amp; Cognitive Language Technology)</i>	2025
NAACL 2025 Registration Grant <i>NAACL 2025 SRW (Student Research Workshop)</i>	2025
Gold Prize (Director's Award of the NIKL) <i>Korean AI Language Proficiency Challenge, NIKL (National Institute of Korean Language)</i>	2024
Excellent Paper Award <i>HCLT 2024 (The 36th Annual Conference on Human &amp; Cognitive Language Technology)</i>	2024
Scholarship for Outstanding Graduate Students <i>POSTECH</i>	2024
Top Engineering Student Award <i>Inha University</i>	2024
Research Scholarship for Undergraduate Researchers <i>Inha University</i>	2023

## TECHNICAL SKILLS

Professional working proficiency Python, PyTorch, transformers, vLLM, Git
Limited working proficiency ADK, CrewAI, Shell Script, Keras, $\text{\LaTeX}$
Elementary proficiency DeepSpeed, TensorFlow, C++, C, MySQL