

Hwaran Lee

CONTACT INFORMATION	NAVER AI Lab at NAVER Cloud	hwaran.lee@gmail.com	hwaranlee.github.io
RESEARCH INTERESTS	My current primary research interests are Ethics, Safety, and Trustworthiness of Large Language Models. I am also interested in controllable language generation, dialog systems, and machine learning for language models.		
EDUCATION	KAIST		Daejeon, South Korea
	Ph.D., Electrical Engineering		Mar. 2013 – Aug. 2018
	Dissertation: <i>Neural Representations for Speech Recognition and Natural Language Generation</i>		
	Advisor: Prof. Soo-Young Lee		
	KAIST		Daejeon, South Korea
	B.S., Mathematical Science		Feb. 2008 – Aug. 2012
	Minor: Financial Engineering Program		
	<i>Magna Cum Laude</i>		
RESEARCH AND WORK EXPERIENCES	NAVER Cloud		Seongnam, South Korea
	<ul style="list-style-type: none">• <i>Leader</i>, Language Research Team, NAVER AI Lab• <i>Tech Leader</i>, NAVER AI Lab• <i>Research Scientist</i>, NAVER AI Lab		Apr. 2023 – Present Jan. 2023 – Mar. 2023 Jan. 2023 – Present
	NAVER		Seongnam, South Korea
	<ul style="list-style-type: none">• <i>Tech Leader</i>, NAVER AI Lab• <i>Research Scientist</i>, NAVER AI Lab		Jul. 2022 – Dec. 2022 Mar. 2021 – Dec. 2022
	SK Telecom		Seoul, South Korea
	<ul style="list-style-type: none">• <i>Research Scientist</i>, T-Brain, AI Center		Nov. 2018 – Feb. 2021
	Brain Science Research Center		Daejeon, South Korea
	<ul style="list-style-type: none">• <i>Undergraduate Researcher</i>		Sep. 2012 – Feb. 2013
PUBLICATIONS	International Journal		
	[J1] Hwaran Lee , Seokhwan Jo, HyungJun Kim, Sangkeun Jung, and Tae-Yoon Kim, “SUMBT+LaRL: Effective Multi-domain End-to-end Neural Task-oriented Dialog System”, <i>IEEE Access</i> , 9 (2021): 116133-116146.		
	[J2] Geonmin Kim, Hwaran Lee , Bo-Kyeong Kim, Sang-Hoon Oh, and Soo-Young Lee, “Unpaired Speech Enhancement by Acoustic and Adversarial Supervision for Speech Recognition”, <i>IEEE Signal Processing Letters</i> , (2019): 159-163.		
	[J3] Ho-Gyeong Kim, Hwaran Lee , Geonmin Kim, Sang-Hoon Oh, and Soo-Young Lee, “Rescoring of N-best Hypotheses using Top-down Selective Attention for Automatic Speech Recognition”, <i>IEEE Signal Processing Letters</i> , (2018): 199-203.		
	[J4] Hwaran Lee , Geonmin Kim, Ho-Gyeong Kim, Sang-Hoon Oh, and Soo-Young Lee, “Deep CNNs Along the Time Axis With Intermap Pooling for Robustness to Spectral Variations”, <i>IEEE Signal Processing Letters</i> 23.10 (2016): 1310-1314.		
	[J5] Hwaran Lee , Nadeem Iqbal, Wonil Chang, and Soo-Young Lee, “A Calibration Method for Eye-Gaze Estimation Systems Based on 3D Geometrical Optics”, <i>IEEE Sensors Journal</i> 13, no. 9 (2013): 3219-3225.		

- [J6] Nadeem Iqbal, **Hwaran Lee**, and Soo-Young Lee, “Smart User Interface for Mobile Consumer Devices Using Model-Based Eye-Gaze Estimation”, *IEEE Transactions on Consumer Electronics* 59, no. 1 (2013): 161-166.

International Conference

- [C1] **Hwaran Lee**^{*}, Seokhee Hong^{*}, Joonsuk Park, Takyoun Kim, Meeyoung Cha, Yejin Choi, Byoungpil Kim, Gunhee Kim, Eun-Ju Lee, Yong Lim, Alice Oh, Sangchul Park, and Jung-Woo Ha, “SQuARe: A Large-Scale Dataset of Sensitive Questions and Acceptable Responses Created through Human-Machine Collaboration”, *ACL*, 2023 (Oral)
- [C2] **Hwaran Lee**^{*}, Seokhee Hong^{*}, Joonsuk Park, Takyoun Kim, Gunhee Kim, and Jung-Woo Ha, “KoSBI: A Dataset for Mitigating Social Bias Risks Towards Safer Large Language Model Applications”, *ACL*, 2023
- [C3] Deokjae Lee, JunYeong Lee, Jung-Woo Ha, Jin-Hwa Kim, Sang-Woo Lee, **Hwaran Lee**, and Hyun Oh Song, “Query-Efficient Black-Box Red Teaming via Bayesian Optimization”, *ACL*, 2023
- [C4] Minbeom Kim, Hwanhee Lee, Kang Min Yoo, Joonsuk Park, **Hwaran Lee**[†], Kyomin Jung[†], “Critic-Guided Decoding for Controlled Text Generation”, *ACL (Findings)*, 2023
- [C5] Miyoung Ko, Ingyu Seong, **Hwaran Lee**, Joonsuk Park, Minsuk Chang, Minjoon Seo, “Beyond Fact Verification: Comparing and Contrasting Claims on Contentious Topics”, *ACL (Findings)*, 2023
- [C6] Hwanhee Lee, Kang Min Yoo, Joonsuk Park, **Hwaran Lee**[†], Kyomin Jung[†], “Masked Summarization to Generate Factually Inconsistent Summaries for Improved Factual Consistency Checking”, *In Findings of the Association for Computational Linguistics: NAACL*, 2022.
- [C7] Kyungjae Lee, Wookje Han, Seung-won Hwang, **Hwaran Lee**, Joonsuk Park, Sang-Woo Lee, “Plug-and-Play Adaptation for Continuously-updated QA”, *In Findings of the Association for Computational Linguistics: ACL*, 2022.
- [C8] John Yoon Young Chung, Wooseok Kim, Kang Min Yoo, **Hwaran Lee**, Eytan Adar, Minsuk Chang, “TaleBrush: Sketching Stories with Generative Pretrained Language Models”, *In Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, 2022.
- [C9] Gi-Cheon Kang, Junseok Park, **Hwaran Lee**, Byoung-Tak Zhang, and Jin-Hwa Kim, “Reasoning Visual Dialog with Sparse Graph Learning and Knowledge Transfer”, *In Findings of the Association for Computational Linguistics: EMNLP*, 2021.
- [C10] **Hwaran Lee**^{*}, Jinsik Lee^{*}, and Tae-Yoon Kim, “SUMBT: Slot-Utterance Matching for Universal and Scalable Belief Tracker”, *The 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- [C11] Geonmin Kim, **Hwaran Lee**, Bo-Kyeong Kim, and Soo-Young Lee, “Compositional Sentence Representation from Character within Large Context Text”, *International Conference on Neural Information Processing (ICONIP)*, 2017.
- [C12] Ho-Gyeong Kim, Jihyeon Roh, **Hwaran Lee**, Geonmin Kim, and Soo-Young Lee, “Active Learning for Large-scale Object Classification: from Exploration to Exploitation” *In Proceedings of the 3rd International Conference on Human-Agent Interaction, (HAI)*, 2015.

[†]corresponding authors

^{*}these authors contributed equally to this work

- [C13] Bo-Kyeong Kim, **Hwaran Lee**, Jihyeon Roh, and Soo-Young Lee, “Hierarchical committee of deep CNNs with exponentially-weighted decision fusion for static facial expression recognition”, *In Proceedings of the 2015 ACM on International Conference on Multimodal Interaction (ICMI)*, 2015.

International Workshop

- [W1] Jaimeen Ahn, **Hwaran Lee**, Jin-Hwa Kim, Alice Oh, “Why Knowledge Distillation Amplifies Gender Bias and How to Mitigate from the Perspective of DistilBERT”, *In Proceedings of the 4rd Workshop on Gender Bias in Natural Language Processing*, 2022.
- [W2] John Yoon Young Chung, Wooseok Kim, Kang Min Yoo, **Hwaran Lee**, Eytan Adar, Minsuk Chang, “TaleBrush: Visual Sketching of Story Generation with Pretrained Language Models”, *CHI EA 22: CHI Conference on Human Factors in Computing Systems Extended Abstracts*, 2022.
- [W3] Geonmin Kim*, **Hwaran Lee***, CheongAn Lee, Eunmi Hong, Byungeun Kim, Soo-Young Lee, “A Deep Chatbot for QA and Chitchat.”, *The Conversational Intelligence Challenge section on NIPS 2017 Competition Track Workshop*, 2017.
- [W4] **Hwaran Lee**, Geonmin Kim, Jihyeon Roh, and Soo-Young Lee, “Learning Tonotopically Organized Auditory Feature-map from Speech by an Intermap Pooling Layer in a Deep CNN”, *15th China-Japan-Korea Joint Workshop on Neurobiology and Neuroinformatics (NBNI)*, 2015. (only abstract)
- [W5] Geonmin Kim, **Hwaran Lee**, Jaemyung Yu, and Soo-Young Lee, “Spoken Sentence Embedding from Character by Jointly Learning Character-level Compositional Word Model and RNN Sentence Encoder”, *15th China-Japan-Korea Joint Workshop on Neurobiology and Neuroinformatics (NBNI)*, 2015. (only abstract)

Pre-prints

- [A1] Taehyun Lee, Seokhee Hong, Jaewoo Ahn, Ilgee Hong, **Hwaran Lee**, Sangdoon Yun, Jamin Shin, Gunhee Kim, “Who Wrote this Code? Watermarking for Code Generation”, arXiv preprint arXiv:2305.15060 (2023), (*Submitted to EMNLP 2023*)
- [A2] Siwon Kim, Sangdoon Yun, **Hwaran Lee**, Martin Gubri, Sungroh Yoon, Seong Joon Oh, “ProPILE: Probing Privacy Leakage in Large Language Models”, (*Submitted to NeurIPS 2023*)

PATENTS

- [P1] Tae-Yoon Kim, Jin Kim, Hyungjoon Kim, Jinsik Lee, **Hwaran Lee**, Heewon Jeon, Seokhwan Jo, “Method and Apparatus for Providing Hybrid Intelligent Customer Consultation”, Korea Patent Application 10-2019-0136035, filed October 2019, Patent Pending.
- [P2] **Hwaran Lee**, Jinsik Lee, Tae-Yoon Kim, “Method and Apparatus for Dialogue State Tracking for Use in Goal-oriented Dialog System”,
- Korea Patent Application 10-2019-0086380, filed July 17, 2019, Patent Pending.
 - PCT/KR2020/008832, filed July 7, 2020, Patent Pending.
 - China Patent 202080051265.8, filed July 7, 2020, Patent Pending.
 - US Patent 17/619,568, filed July 7, 2020, Patent Pending.

HONORS AND AWARDS	<ul style="list-style-type: none"> • Annual Roll Award, KAIST EE Apr. 2018 • Ranked 3rd, ConvAI challenge, NIPS 2017 Competition Track Workshop 2017 • Challenge Winner, ICMI EmotiW2015 2015 • Best Paper Award, HAI 2015 • Qualcomm Innovation Award 2015 • BK21 Plus Financial Support for Graduates Long Term Training May. 2014 • KAIST Graduate Scholarship Mar. 2013 - Aug. 2018 • Australian Endeavour Student Exchange Grant (AUD\$ 5000), The University of Queensland Apr. 2011 • National Excellence Scholarship, KOSAF Feb.2008 - Feb. 2012
ACADEMIC SERVICES	<ul style="list-style-type: none"> • Organizing Committee <ul style="list-style-type: none"> • ACM FAccT 2022 CRAFT HyperscaleFAccT • Area Chair <ul style="list-style-type: none"> • NeurIPS 2023 Datasets & Benchmarks • Reviewer <ul style="list-style-type: none"> • ARR 2021-2022, ACL 2021-2023, EMNLP 2021-2023, COLING 2020, 2022 • NeurIPS 2021-2023, ICLR 2021-2022 • WWW 2022 • ACL'22 In2Writing Workshop • Samsung Humantech Paper Awards Committee 2020 • Qualcomm Innovation Awards Committee 2019 • Speech Communication 2019 • IEEE Transactions on Neural Networks and Learning Systems 2017 - 2018 • Neural Processing Letters 2015
OUTSIDE ACTIVITIES	<ul style="list-style-type: none"> • Committee member of the 2nd Forum on Artificial Intelligence Ethics and Policy, organized by the Ministry of Science and ICT, South Korea. 2023 • Organizing committee of AI Ethics Forum for Human at NAVER 2022
INVITED TALKS	<ul style="list-style-type: none"> • Ethical Problems in Language Models <ul style="list-style-type: none"> Intellectual Property High Court, Daejeon, Apr. 2023 GIST, Gwang-ju, Feb. 2023 Hanyang Univ., Seoul, Sep. 2022 KAIST, Daejeon, Jun. 2022 SNU, Seoul, Jun. 2022 • Introduction to deep learning for dialogue systems <ul style="list-style-type: none"> Inha Univ., Seoul, Oct. 2020 Yeonsei Univ., Seoul, Jun. 2020 Sookmyung Women's Univ., Seoul, Oct. 2019 • Toward end-to-end neural dialog systems for multi-domain task completion <ul style="list-style-type: none"> KAIST, Daejeon, Dec. 2019
RESEARCH PROJECTS (BEFORE 2021)	<p>Meta Learner Project 2020 - 2021</p> <ul style="list-style-type: none"> • Funded by SK Telecom <ul style="list-style-type: none"> – Researched and developed semi-supervised learning algorithms and modules for AutoTrainer – Tech leader for research and development of AutoML for various natural language processing tasks <p>End-to-end Goal-Oriented Dialog Systems 2018 - 2019</p> <ul style="list-style-type: none"> • Funded by SK Telecom

- Researched and developed a scalable multi-domain natural language understanding
- Researched and developed end-to-end neural goal-oriented dialog systems
- Participated in the End-to-End Multi-Domain Task-Completion Task (DSTC8 Track1 Task1 at AAI 2020) and ranked 6th place

A Deep Chatbot for QA and Chitchat 2017

- Project Leader, Funded by *Institute for Information & Communications Technology Promotion (IITP)*
 - Researched and developed an article-based chatbot that carry on both question-answering and chitchat conversations
 - **Ranked 3rd Place** in the Conversational Intelligence Challenge of NIPS 2017 Live Competition

Deep Learning Based Korean Understanding and Applications 2015 – 2018

- Project Leader, Funded by HANCOM Inc.
 - Researched and developed Korean morpheme- and syllable-level language models, sentence and document models
 - Researched and developed continual learning algorithms for language models
 - KoreanLM-demo (github)

Speech Feature Extraction Based on Deep Learning for Continuous Speech Recognition Systems 2013 – 2014

- Project Leader, Funded by LG Electronics
 - Developed convolutional neural networks for acoustic models of continuous English speech
 - Programmed C/C++ and CUDA based on KALDI toolkit
 - KALDI-CNN(github)

Development of Dialog-based Spontaneous Speech Interface Technology on Mobile Platforms 2013 – 2014

- Funded by Electronics and Telecommunications Research Institute (*ETRI*)
 - Developed deep belief networks and auto-encoders for acoustic models of Korean syllable data
 - Programmed models in C/C++ and CUDA languages

TEACHING
EXPERIENCE
(BEFORE 2017)

Teaching Assistant

- ChatBot PyTorch Implementation, KB-KAIST AI Sep. 2017
- EE476 Audio-Visual Perception Models Spring 2016, 2017
- EE538 Neural Networks Fall 2015, 2016, 2017

Graduate Mentor

- Undergraduate Mentoring Program, KAIST Counselling Center
Fall 2014, Spring 2016