

Seanie Lee

DATE OF BIRTH: 1992.04.17, NATIONALITY: KOREAN

126 Yangjae-dong, Seocho District, Seoul

□ (+82) 10-4475-2273 | □ lsnfamily02@kaist.ac.kr | □ seanie12.github.io | □ seanie12 | □ Seanie Lee

Education

KAIST (Korea Advanced Institute of Science and Technology)

PH.D IN ARTIFICIAL INTELLIGENCE

Daejeon, S.Korea

Mar. 2022 - present

- Supervised by [Sung Ju Hwang](#) and [Juho Lee](#)

• Research interest: AI safety, responsible AI, and evaluation

KAIST (Korea Advanced Institute of Science and Technology)

M.S. IN ARTIFICIAL INTELLIGENCE

Daejeon, S.Korea

Mar. 2020 - Feb. 2022

- Supervised by [Sung Ju Hwang](#) and [Juho Lee](#)

• Master Thesis: [Data augmentation for natural language processing](#)

Yonsei University

Seoul, S.Korea

B.A. IN LIBRARY AND INFORMATION SCIENCE

Mar. 2011 - Feb. 2018

Experience

Apple

Seattle, US

INTERNSHIP

October 2025 - May 2026

- Machine Learning Research, hosted by [Raviteja Vemulapalli](#).
- Synthetic data generation for tool-calling LLMs.

Krafton

Seoul, Korea

INTERNSHIP

July 2025 - Oct 2025

- Research Internship
- Safety alignment of long reasoning models.

Mila

Montreal, Canada

INTERNSHIP

January 2024 - June 2024

- Research internship at Mila, advised by [Yoshua Bengio](#).
- Robust red-teaming of LLMs.

Apple

Cambridge, UK

INTERNSHIP

May 2023 - September 2023

- Research internship at Siri team, hosted by [Anders Johannsen](#).
- Few-shot example retrieval for ICL.

National University of Singapore (NUS)

Singapore

INTERNSHIP

July 2022 - September 2022

- Remote internship at [Deep Learning](#) lab, supervised by [Kenji Kawaguchi](#).
- Regularization for continual pretraining.

Awards

2023 **Apple AI/ML PhD Fellowship**, Recipient of [Apple Scholars in AI/ML](#)

Cupertino, US

2022 **Google Travel Grant**, NeurIPS 2022

US

2019 **Silver Medal**, Named Entity Recognition in [NAVER NLP Challenge](#)

Seoul, Korea

Presentation

Seminar at Korea University.

Seoul, Korea

PRESENTATION OF LARGE SCALE SET-ENCODING

May. 2025

- Synthetic Data Generation for LLM Safeguards
- ICLR 2025, ACL Findings 2025

Seminar at Hanyang University.

Seoul, Korea

April. 2025

PRESENTATION OF LARGE SCALE SET-ENCODING

- Synthetic Data Generation for LLM Safeguards
- ICLR 2025, ACL Findings 2025

Tech. Talk, Nuremberg Institute of Technology Georg Simon Ohm.

Nürnberg, Germany

Oct. 2023

PRESENTATION OF LARGE SCALE SET-ENCODING

- Scalable Set Encoding with Universal Mini-Batch Consistency and Unbiased Full Set Gradient Approximation
- ICML 2023

Tech. talk, Samgsung SDS.

Seoul, South Korea

22.May. 2023

PRESENTATION OF LARGE SCALE SET-ENCODING

- Scalable Set Encoding with Universal Mini-Batch Consistency and Unbiased Full Set Gradient Approximation
- ICML 2023

Tech. talk, NAVER corp.

Online, South Korea

04.Dec. 2020

PRESENTATION OF INFO-HCVAE

- Generating Diverse and Consistent QA pairs from Contexts with Information-Maximizing Hierarchical Conditional VAEs
- ACL 2020 Long paper

Selected Publication

(* indicates equal contribution)

PREPRINT

HoliSafe: Holistic Safety Benchmarking and Modeling with Safety Meta Token for Vision-Language Model

Arxiv

YOUNGWAN LEE, KANGSAN KIM, KWANYONG PARK, ILCAHE JUNG, SOOJIN JANG, SEANIE LEE, YONG-JU LEE AND SUNG JU

HWANG

- [\[paper\]](#)[\[code\]](#)

CONFERENCES

Learning Diverse Attacks on Large Language Models for Robust Red-teaming and Safety Tuning

ICLR

SEANIE LEE, MINSU KIM, LYNN CHERIF, DAVID DOBRE, JUHO LEE, SUNG JU HWANG, KENJI KAWAGUCHI, GAUTHIER GIDEL, YOSHUA BENGIO, NIKOLAY MALKIN, MOKSH JAIN

2025

- [\[paper\]](#)[\[code\]](#)

HarmAug: Effective Data Augmentation for Knowledge Distillation of Safety Guard Models

ICLR

SEANIE LEE*, HAEBIN SEONG*, DONG BOK LEE, MINKI KANG, XIAOYIN CHEN, DOMINIK WAGNER, YOSHUA BENGIO, JUHO LEE, SUNG JU HWANG

2025

- [\[paper\]](#)[\[code\]](#)

SafeRoute: Adaptive Model Selection for Efficient and Accurate Safety Guardrails in Large Language Models

ACL Findings

SEANIE LEE*, DONG BOK LEE*, DOMINIK WAGNER, MINKI KANG, HAEBIN SEONG, TOBIAS BOCKLET, JUHO LEE, SUNG JU

HWANG

- [\[paper\]](#)[\[code\]](#)

Reliable Decision-Making via Calibration-Oriented Retrieval-Augmented Generation

NeurIPS

CHAEYUN JANG, DEUKHWAN CHO, SEANIE LEE, JUHO LEE AND HYUNGI LEE

2025

- [\[paper\]](#)[\[code\]](#)

FedSVD: Adaptive Orthogonalization for Private Federated Learning with LoRA

NeurIPS

SEANIE LEE*, SANGWOO PARK*, DONG BOK LEE*, DOMINIK WAGNER, HAEBIN SEONG, TOBIAS BOCKLET, JUHO LEE, SUNG JU

HWANG

- [\[paper\]](#)[\[code\]](#)

Publication

(* indicates equal contribution)

PREPRINT

Rethinking Reward Models for Multi-Domain Test-Time Scaling

Arxiv

DONG BOK LEE*, **SEANIE LEE***, SANGWOO PARK, MINKI KANG, JINHEON BAEK, DONGKI KIM, DOMINIK WAGNER, JIONGDAO JIN, HEEJUN LEE, TOBIAS BOCKLET, JINYU WANG, JINGJING FU, SUNG JU HWANG, JIANG BIAN AND LEI SONG

2025

- [paper][code]

HoliSafe: Holistic Safety Benchmarking and Modeling with Safety Meta Token for Vision-Language Model

Arxiv

YOUNGWAN LEE, KANGSAN KIM, KWANYONG PARK, ILCAHE JUNG, SOOJIN JANG, **SEANIE LEE**, YONG-JU LEE AND SUNG JU HWANG

2025

- [paper][code]

CONFERENCES

FedSVD: Adaptive Orthogonalization for Private Federated Learning with LoRA

NeurIPS

SEANIE LEE*, SANGWOO PARK*, DONG BOK LEE*, DOMINIK WAGNER, HAEBIN SEONG, TOBIAS BOCKLET, JUHO LEE, SUNG JU HWANG

2025

- [paper][code]

Distilling LLM Agent into Small Models with Retrieval and Code Tools

NeurIPS Spotlight

MINKI KANG, JONGWON JEONG, **SEANIE LEE**, JAEWOONG CHO AND SUNG JU HWANG

2025

- [paper][code]

Reliable Decision-Making via Calibration-Oriented Retrieval-Augmented Generation

NeurIPS

CHAEYUN JANG, DEUKHWAN CHO, **SEANIE LEE**, JUHO LEE AND HYUNGI LEE

2025

- [paper][code]

Trajectory Balance with Asynchrony: Decoupling Exploration and Learning for Fast, Scalable LLM Post-Training

NeurIPS

BRIAN R. BARTOLDSON, SIDDARTH VENKATRAMAN, JAMES DIFFENDERFER, MOKSH JAIN, TAL BEN-NUN, **SEANIE LEE**, MINSU KIM, JOHAN OBANDO-CERON, YOSHUA BENGIO AND BHAVYA KAILKHURA

2025

- [paper][code]

SafeRoute: Adaptive Model Selection for Efficient and Accurate Safety Guardrails in Large Language Models

ACL Findings

SEANIE LEE*, DONG BOK LEE*, DOMINIK WAGNER, MINKI KANG, HAEBIN SEONG, TOBIAS BOCKLET, JUHO LEE, SUNG JU HWANG

2025

- [paper][code]

Personalized Fine-Tuning with Controllable Synthetic Speech from LLM-Generated Transcripts for Dysarthric Speech Recognition

Interspeech

DOMINIK WAGNER, ILJA BAUMANN, NATALIE ENGERT, **SEANIE LEE**, ELMAR NÖTH, KORBINIAN RIEDHAMMER AND TOBIAS BOCKLET

2025

- [paper]

HarmAug: Effective Data Augmentation for Knowledge Distillation of Safety Guard Models

ICLR

SEANIE LEE*, HAEBIN SEONG*, DONG BOK LEE, MINKI KANG, XIAOYIN CHEN, DOMINIK WAGNER, YOSHUA BENGIO, JUHO LEE, SUNG JU HWANG

2025

- [paper][code]

Learning Diverse Attacks on Large Language Models for Robust Red-teaming and Safety Tuning

ICLR

SEANIE LEE, MINSU KIM, LYNN CHERIF, DAVID DOBRE, JUHO LEE, SUNG JU HWANG, KENJI KAWAGUCHI, GAUTHIER GIDEL, YOSHUA BENGIO, NIKOLAY MALKIN, MOKSH JAIN

2025

- [\[paper\]](#)[\[code\]](#)

Optimized Speculative Sampling for GPU Hardware Accelerators

EMNLP

DOMINIK WAGNER, SEANIE LEE, ILJA BAUMANN, PHILIPP SEEGER, KORBINIAN RIEDHAMMER, TOBIAS BOCKLET

2024

- [\[paper\]](#)[\[code\]](#)

Drug Discovery with Dynamic Goal-aware Fragment

ICML

SEUL LEE, SEANIE LEE, KENJI KAWAGUCHI, SUNG JU HWANG

2024

- [\[paper\]](#)[\[code\]](#)

Effective and Efficient Conversation Retrieval for Dialogue State Tracking with Implicit Text Summaries

NAACL

SEANIE LEE, JIANPENG CHENG, JORIS DRIESSEN, ALEXANDRU COCA, ANDERS JOHANNSEN

2024

- [\[paper\]](#)

Self-Supervised Dataset Distillation for Transfer Learning

ICLR

DONG BOK LEE*, SEANIE LEE*, JOONHO KO, KENJI KAWAGUCHI, JUHO LEE, SUNG JU HWANG

2024

- [\[paper\]](#)[\[code\]](#)

DiffusionNAG: Task-guided Neural Architecture Generation with Diffusion Models

ICLR

SOHYUN AHN*, HAYEON LEE*, JAEHYEONG JO, SEANIE LEE, SUNG JU HWANG

2024

- [\[paper\]](#)[\[code\]](#)

Scalable Set Encoding with Universal Mini-Batch Consistency and Unbiased Full Set Gradient Approximation

ICML

JEFFREY WILLETT*, SEANIE LEE*, BRUNO ANDREIS, KENJI KAWAGUCHI, JUHO LEE, SUNG JU HWANG

2023

- [\[paper\]](#)[\[code\]](#)

Margin-based Neural Network Watermarking

ICML

BYUNGJOO KIM, SUYOUNG LEE, SEANIE LEE, SOOEL SON, SUNG JU HWANG

2023

- [\[paper\]](#)

Self-Supervised Set Representation Learning for Unsupervised Meta-Learning

ICLR

DONG BOK LEE*, SEANIE LEE*, KENJI KAWAGUCHI, YUNJI KIM, JIHWA BANG, JUNG-WOO HA, SUNG JU HWANG

2023

- [\[paper\]](#)

Self-Distillation for Further Pre-training of Transformers

ICLR

SEANIE LEE, MINKI KANG, JUHO LEE, SUNG JU HWANG, KENJI KAWAGUCHI

2023

- [\[paper\]](#)[\[code\]](#)

Set-based Meta-Interpolation for Few-Task Meta-Learning

NeurIPS

SEANIE LEE*, BRUNO ANDREIS*, KENJI KAWAGUCHI, SUNG JU HWANG

2022

- [\[paper\]](#) [\[code\]](#)

On Divergence Measures for Bayesian Pseudocoresets

NeurIPS

BALHAE KIM, JUNGWON CHOI, SEANIE LEE, YOONHO LEE, JUNG-WOO HA, JUHO LEE

2022

- [\[paper\]](#)

Set Based Stochastic Subsampling

ICML

BRUNO ANDREIS, SEANIE LEE, A. TUAN NGUYEN, JUHO LEE, EUNHO YANG, SUNG JU HWANG

2022

- [\[paper\]](#)

Sequential Reptile: Inter-Task Gradient Alignment for Multilingual Learning

ICLR

SEANIE LEE*, HAE BEOM LEE*, JUHO LEE, SUNG JU HWANG

2022

- [\[paper\]](#)

Learning to Perturb Word Embeddings for Out-of-distribution QA

ACL

SEANIE LEE*, MINKI KANG*, JUHO LEE, SUNG JU HWANG

2021

- [\[paper\]](#)[\[code\]](#)

Contrastive Learning with Adversarial Perturbations for Conditional Text Generation

ICLR

SEANIE LEE*, DONG BOK LEE*, SUNG JU HWANG

2021

- [\[paper\]](#)[\[code\]](#)

Meta-GMVAE: Mixture of Gaussian VAE for Unsupervised Meta-Learning

DONG BOK LEE, DONGCHAN MIN, **SEANIE LEE**, SUNG JU HWANG

- [\[paper\]](#)[\[code\]](#)

ICLR

Generating Diverse and Consistent QA pairs from Contexts with Information-Maximizing Hierarchical Conditional VAEs

DONG BOK LEE*, **SEANIE LEE***, WOOTAE JEONG, DONGHWAN KIM, SUNG JU HWANG

- [\[paper\]](#) [\[code\]](#)[\[video\]](#)

g2pM: A Neural Grapheme-to-Phoneme Conversion Package for Mandarin Chinese Based on a New Open Benchmark Dataset

KYUBYONG PARK*, **SEANIE LEE***

- [\[paper\]](#)[\[code\]](#)

INTERSPEECH

References

Sung Ju Hwang

ASSOCIATE PROFESSOR IN KAIST.

e-mail: sjhwang82@kaist.ac.kr.

Advisor

2020-2025

Juho Lee

ASSOCIATE PROFESSOR IN KAIST.

e-mail: juholee@kaist.ac.kr.

Advisor

2020-2025

Yoshua Bengio

FULL PROFESSOR AT UNIVERSITÉ DE MONTRÉAL AND SCIENTIFIC DIRECTOR OF MILA – QUEBEC AI INSTITUTE.

e-mail: yoshua.bengio@mila.quebec.

Collaborator

2024-2025

Kenji Kawaguchi

PRESIDENTIAL YOUNG PROFESSOR IN THE DEPARTMENT OF COMPUTER SCIENCE AT NUS.

e-mail: kenji@comp.nus.edu.sg

Collaborator

2022-present

Nikolay Malkin

CHANCELLOR'S FELLOW AT UNIVERSITY OF EDINBURGH, SCHOOL OF INFORMATICS

e-mail: nmalkin@ed.ac.uk

Collaborator

2024