

Seanie Lee

DATE OF BIRTH: 1992.04.17, NATIONALITY: KOREAN

126 Yangjae-dong, Seocho District, Seoul

□ (+1) 206-468-5851 | ✉ lsfamily02@gmail.com | 🌐 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

Research Interest

- Data-Centric AI [2, 4, 5, 9, 16, 19, 20, 25, 27, 28]
- Responsible and Safe AI [19, 20, 22, 24, 23, 26, 27]
- Efficient ML [7, 13, 16, 18, 20, 22, 25]
- Agentic AI [25]

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

Publication

(* indicates equal contribution)

PREPRINT

[28] 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

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

[27] 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

[26] 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

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

[25] 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

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

[24] Reliable Decision-Making via Calibration-Oriented Retrieval-Augmented Generation

NeurIPS

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

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

[23] 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

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

[22] 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

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

[21] 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

- [\[paper\]](#)

[20] 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

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

[19] 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

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

[18] Optimized Speculative Sampling for GPU Hardware Accelerators

EMNLP

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

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

[17] Drug Discovery with Dynamic Goal-aware Fragment*ICML*

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

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

[16] 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

- [\[paper\]](#)

[15] Self-Supervised Dataset Distillation for Transfer Learning*ICLR*

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

2024

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

[14] DiffusionNAG: Task-guided Neural Architecture Generation with Diffusion Models*ICLR*

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

2024

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

[13] Scalable Set Encoding with Universal Mini-Batch Consistency and Unbiased Full Set Gradient Approximation*ICML*

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

2023

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

[12] Margin-based Neural Network Watermarking*ICML*

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

2023

- [\[paper\]](#)

[11] 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

- [\[paper\]](#)

[10] Self-Distillation for Further Pre-training of Transformers*ICLR*

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

2023

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

[9] Set-based Meta-Interpolation for Few-Task Meta-Learning*NeurIPS*

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

2022

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

[8] On Divergence Measures for Bayesian Pseudocoresets*NeurIPS*

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

2022

- [\[paper\]](#)

[7] Set Based Stochastic Subsampling*ICML*

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

2022

- [\[paper\]](#)

[6] Sequential Reptile: Inter-Task Gradient Alignment for Multilingual Learning*ICLR*

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

2022

- [\[paper\]](#)

[5] Learning to Perturb Word Embeddings for Out-of-distribution QA*ACL*

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

2021

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

[4] Contrastive Learning with Adversarial Perturbations for Conditional Text Generation*ICLR*

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

2021

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

[3] Meta-GMVAE: Mixture of Gaussian VAE for Unsupervised Meta-Learning*ICLR*

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

2021

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

[2] Generating Diverse and Consistent QA pairs from Contexts with Information-Maximizing Hierarchical Conditional VAEs*ACL*

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

2020

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

KYUBYONG PARK*, SEANIE LEE*

2020

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

Publication

(* indicates equal contribution)

PREPRINT

[28] 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

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

[27] 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

[26] 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

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

[25] 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

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

[24] Reliable Decision-Making via Calibration-Oriented Retrieval-Augmented Generation

NeurIPS

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

2025

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

[23] 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

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

[22] 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

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

[21] 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

- [\[paper\]](#)

[20] 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

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

[19] 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]	
[18] Optimized Speculative Sampling for GPU Hardware Accelerators	EMNLP
DOMINIK WAGNER, SEANIE LEE , ILJA BAUMANN, PHILIPP SEEBERGER, KORBINIAN RIEDHAMMER, TOBIAS BOCKLET	2024
• [paper] [code]	
[17] Drug Discovery with Dynamic Goal-aware Fragment	ICML
SEUL LEE, SEANIE LEE , KENJI KAWAGUCHI, SUNG JU HWANG	2024
• [paper] [code]	
[16] 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]	
[15] 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]	
[14] DiffusionNAG: Task-guided Neural Architecture Generation with Diffusion Models	ICLR
SOHYUN AHN*, HAYEON LEE*, JAEHYEONG JO, SEANIE LEE , SUNG JU HWANG	2024
• [paper] [code]	
[13] 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]	
[12] Margin-based Neural Network Watermarking	ICML
BYUNGJOO KIM, SUYOUNG LEE, SEANIE LEE , SOOEL SON, SUNG JU HWANG	2023
• [paper]	
[11] 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]	
[10] Self-Distillation for Further Pre-training of Transformers	ICLR
SEANIE LEE , MINKI KANG, JUHO LEE, SUNG JU HWANG, KENJI KAWAGUCHI	2023
• [paper] [code]	
[9] Set-based Meta-Interpolation for Few-Task Meta-Learning	NeurIPS
SEANIE LEE *, BRUNO ANDREIS*, KENJI KAWAGUCHI, SUNG JU HWANG	2022
• [paper] [code]	
[8] On Divergence Measures for Bayesian Pseudocoresets	NeurIPS
BALHAE KIM, JUNGWON CHOI, SEANIE LEE , YOONHO LEE, JUNG-WOO HA, JUHO LEE	2022
• [paper]	
[7] Set Based Stochastic Subsampling	ICML
BRUNO ANDREIS, SEANIE LEE , A. TUAN NGUYEN, JUHO LEE, EUNHO YANG, SUNG JU HWANG	2022
• [paper]	
[6] Sequential Reptile: Inter-Task Gradient Alignment for Multilingual Learning	ICLR
SEANIE LEE *, HAE BEOM LEE*, JUHO LEE, SUNG JU HWANG	2022
• [paper]	
[5] Learning to Perturb Word Embeddings for Out-of-distribution QA	ACL
SEANIE LEE *, MINKI KANG*, JUHO LEE, SUNG JU HWANG	2021
• [paper] [code]	
[4] Contrastive Learning with Adversarial Perturbations for Conditional Text Generation	ICLR
SEANIE LEE *, DONG BOK LEE*, SUNG JU HWANG	2021
• [paper] [code]	

[3] Meta-GMVAE: Mixture of Gaussian VAE for Unsupervised Meta-Learning

ICLR

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

2021

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

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

ACL

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

2020

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

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

INTERSPEECH

KYUBYONG PARK*, SEANIE LEE*

2020

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

References

Sung Ju Hwang

Advisor

ASSOCIATE PROFESSOR IN KAIST.

2020-2025

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

Juho Lee

Advisor

ASSOCIATE PROFESSOR IN KAIST.

2020-2025

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

Yoshua Bengio

Collaborator

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

2024-2025

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

Kenji Kawaguchi

Collaborator

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

2022-present

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

Nikolay Malkin

Collaborator

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

2024

e-mail: nmalkin@ed.ac.uk