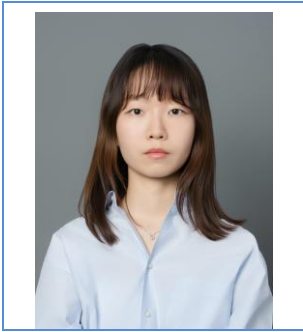


딥러닝 자연어 Researcher / Engineer



이름	김도연
생년월일	1997년 8월 13일
핸드폰	010-2388-7269
주소	경기도 성남시 성남대로 51
이메일	dyeon.kim@snu.ac.kr

학력사항

재학기간	학력사항	전공	학점
2013. 03 - 2016. 02	낙생고등학교	인문	-
2016. 03 - 2020. 02	한양대학교 ERICA 캠퍼스	전자공학부	4.16 / 4.5
2020. 03 - 2022. 08	광주과학기술원 (GIST)	AI 대학원	3.78 / 4.5
2026.03 -	서울대학교	인공지능 협동과정	-

Research interests 및 Topic

기후 Foundation model 연구: Transformer, Multi-modal

주요 경력

기간	구분	기관명
2023.03 - 2025.07	경력사항	주식회사 파수 (Fasoo) 개발센터 AI 팀
		구축형 (Enterprise LLM) ELLM 프로그램 개발자

참여 프로젝트

직책	전임	프로젝트 기간	2024.07 - 2025.07
프로젝트 명		FRB (Fasoo RAG Benchmark) LLM 평가 시스템	
업무 내용		ELLM studio 의 프로젝트 중 하나인 FRB 의 전체적인 기획을 진행하였고, 파이썬 서버를 구현하였다. <ul style="list-style-type: none"> - 기존 ELLM 웹페이지에 FRB 기능을 추가하기 위해 와이어프레임을 작성하고 Use Case 를 작성. - 개발자로서 ragas 라이브러리에 Unanswerable 과 Format evaluation 등의 새 	

	<p>로운 평가지표를 추가하였고, 한국어 지원이 가능하도록 프롬프트를 만들고 refactoring.</p> <ul style="list-style-type: none"> - 파이썬 서버 개발자로서 Flask 를 사용하여 웹서버와 연동. <p>주요 기술 스택: python, Flask</p>
업무 성과	<p>Ellm studio 의 주요기능 중 하나로서 2025.04.22 에 Fasoo Digital Intelligence(FDI) 에 version 1 을 시연하였다.</p> <p>기여도: 25% (UI 개발 1 명, 디자이너 1aud, 웹서버 개발 1 명, 파이썬 서버 개발 및 기획자 1 명)</p>

직책	전임	프로젝트 기간	2024.07 – 2024.10
프로젝트 명		도메인 특화 생성 및 검색 기능 평가 시스템	
업무 내용		<p>Ellm v1.2: 모델 성능 평가를 위한 evaluation 데이터셋을 결정하고 benchmark 자동화 프로그램을 설계하고 구현하였다.</p> <ul style="list-style-type: none"> - 기존 모델을 평가하기 위해 사용하는 evaluation benchmark (SimpleEval 등) 와 RAG 를 평가하는 ragas 적용. - 새로운 평가지표 Unanswerable, Format evaluation 추가. - 도메인 fine-tuning 을 평가하기 위해 기존 평가 데이터 생성 방법인 RGB 와 ragas 를 도메인 평가에 응용하여 적용. <p>주요 기술 스택: Python, Pytorch, vllm, FastChat 사용</p>	
업무 성과		기여도: 100% (개발 1 명)	

직책	사원 3 년차	프로젝트 기간	2024.05 – 2024.06
프로젝트 명		Ellm v1.1 한국어 데이터셋 유해표현 필터링	
업무 내용		<p>Ellmv1.1 의 한국어 학습을 위한 Korean mC4 데이터셋을 hate speech filtering 과 bad word filtering 을 이용해 유해 표현 처리하였다.</p> <p>주요 기술 스택: Pytorch, re</p>	
업무 성과		<p>약 16K mC4 학습데이터셋 기준 20.1%의 욕설 및 혐오 데이터를 필터링할 수 있었다. 이후 한국어 학습 데이터셋 전처리에 유해표현 처리를 추가하여 모든 데이터셋에 적용하였다.</p> <p>기여도 : 100% (개발 1 명)</p>	

직책	사원 3 년차	프로젝트 기간	2024.03 – 진행중
프로젝트 명		Ellm v1.1 Embedding 한국어 학습 방법 연구	
업무 내용		<p>기존 영어 기반의 모델에서 한국어 임베딩 레이어와 토큰이 효율적이지 않게 학습됨을 지적하고 이를 해결하기 위해 한국어 토큰을 새로 추가하고 임베딩 레이어를 일부 학습하는 방법 제시 및 학습 실험.</p> <p>주요 기술 스택: Pytorch, transformers</p>	

업무 성과	<p>이 후, ELLM v1.4 에서 논문 "Efficient and Effective Vocabulary Expansion Towards Multilingual Large Language Models" 기술을 일부 참조하여 한국어 학습에 적용하는 중이다.</p> <p>영어 샘플에 대해서는 추론결과가 동일하고 한국어 샘플에 대해서 추론결과가 개선됨을 확인하였다.</p> <p>기여도 : 50% (개발 2 명)</p>
-------	---

직책	사원 3 년차	프로젝트 기간	2023.08 – 2024.04
프로젝트 명		ELLM v1 개발	
업무 내용	<p>ELLM v1 의 요약 기능 구현을 위한 LLM 기술(Long Context, PEFT 학습(LoRA, Prompt Tuning, IA3))을 리뷰하고 성능 비교 실험을 진행하여 한국어 학습을 진행하였다.</p> <p>이 후 ELLM v1 prompt engineering</p> <ul style="list-style-type: none"> - version 1 기준으로 Agent 구현을 위한 task 에 대한 결과물 생성을 위해 모델의 system 및 task-specific instruction (프롬프트)을 튜닝하여 ELLM v1 프로그램에 적용하고 출시 <p>주요 기술 스택: Pytorch, transformers, peft</p>		
업무 성과	<p>이후 ELLM 에 사용되는 모든 모델 학습을 LoRA 적용</p> <p>Task classification 기준 정확도 58.3%에서 91.6%까지 증가.</p> <p>기여도: 16.7% (개발 6 명)</p>		

직책	인턴	프로젝트 기간	2023.04 -2023.07 (3 개월)
프로젝트 명		NER 모델에 Label Increment Learning 적용	
업무 내용	<p>FasooAIR 의 NER task classification 모델에서 continual learning (label increment) 실험 및 구현</p> <ul style="list-style-type: none"> - 이미 학습된 기존 라벨들의 성능을 떨어트리지 않고 새로운 라벨을 학습할 수 있도록 하는 Label increment learning 을 진행 - EWC(Elastic Weight Consolidation) 과 리허설(Rehearsal)을 적용하여 label 을 추가학습 <p>주요 기술 스택: Pytorch</p>		
업무 성과	<p>리허설과 ewc 를 적용하였을 때 F1 score 기준 기존 모델의 라벨 성능을 평균 점수 기준 0.93 을 유지하면서 추가한 라벨의 f1 점수를 0.96 까지 올릴 수 있었다.</p> <p>기여도: 100% (개발 1 명)</p>		

Skills

Programming	Library
Python, Pytorch	Transformers, PEFT, Flask

어학능력

시험 종류	점수
TOEIC	880

Development	Collaboration
Git, VS code, Linux	Github, Notion, Figma

기타 사항

연구

지도 교수	김강일 (IRR Lab)	연구 기간	2023.08 – 2023.04
논문 제목		졸업논문: Vector-Quantization for representation on Transformer	
연구 내용		<p>본 논문은 Low resource 데이터셋 환경에서 Transformer 모델 학습을 위해 SM-VQ(Soft-Mixed Vector Quantization) 모듈을 이용해 모델의 representation 일반화 연구하였다.</p> <ul style="list-style-type: none"> - Transformer 에 맞는 새로운 모듈인 SM-VQ(Soft-Mixed Vector Quantization) 을 제안. 	
연구 성과		<p>결과: Attention score 에 적용한 결과 low-resource 환경에서 일반 트랜스포머보다 validation 정확도가 약 5% 높아진 것을 확인하여 작은 데이터셋에서 일반화 효과를 얻는 것을 확인할 수 있었다. 또한 random variable 에 대한 안정성이 매우 높아진 것을 확인할 수 있었다.</p>	

지도 교수	김강일 (IRR Lab)	연구 기간	2023.08 – 2023.04
논문 제목		2 저자 논문: Feature Structure Distillation with Centered Kernel Alignment in BERT Transferring	
연구 내용		<p>본 논문은 Teacher 지식에 대한 부정확한 학습을 줄이기 위해 student 모델의 representation 의 flexibility 를 줄이는 새로운 Knowledge Distillation 방법 제안한다.</p> <ul style="list-style-type: none"> - 그 중 Full-batch relation 에 대한 global inter-feature 정보를 클러스터링 알고리즘을 이용하여 distillation 하는 방법을 고안. 	
연구 성과		<p>결과: BERT 모델을 이용하여 GLUE dataset 의 7 개 데이터셋에서 state-of-the-art 성능을 확인하였다. ESA 저널 논문 게재 및 특허 취득</p> <p>Hee-Jun Jung, Doyeon Kim, Seung-Hoon Na, Kangil Kim, Feature structure distillation with Centered Kernel Alignment in BERT transferring, Expert Systems with Applications, Volume 234, 2023, 120980, ISSN 0957-4174, https://doi.org/10.1016/j.eswa.2023.120980.</p>	

기타

주최	Junction Asia 2022	프로젝트 기간	2022.08.19 – 2022.08.21
프로젝트 제목		Microsoft Teams 를 이용한 채팅내용 QA 요약 챗봇 Sprout	

프로젝트 내용	<p>Microsoft 트랙: Azure, Microsoft Teams, Power Platform 을 활용한 협업 앱 개발 역할: 팀장, AI 엔지니어링</p> <p>Microsoft Teams 프로그램 내 채팅 요약 및 질의응답 기능 구현</p> <ul style="list-style-type: none"> - Huggingface API 를 활용하여 SQuAD 데이터셋으로 파인튜닝된 Roberta 모델을 설정하여 채팅 내용 기반 질문 응답 기능 구현 - NLPCloud API 를 활용하여 BART 모델을 설정하고 채팅 내용 요약 기능 구현
프로젝트 성과	Github: 🌱 Sprout-github

주최	광주과학기술원	프로젝트 기간	2021.03 - 2021.08
프로젝트 제목	릴레이 소설 생성 프로그램		
프로젝트 내용	<p>사용자가 딥러닝 모델과 번갈아 가며 소설을 함께 창작할 수 있도록 하는 프로그램 제작하였다.</p> <ul style="list-style-type: none"> - 저작권이 완료된 한국어 소설 4.2MB 을 bs 라이브러리를 이용하여 웹크롤링 및 전처리 - KoGPT-2 를 1 단락(2~3 문장) 단위로 특수 토큰을 추가하여 fine tuning. - 홈페이지는 파이썬 Flask 와 html 을 이용 <p>사용 기술: Pytorch, Python Beautiful Soup(웹 크롤링), Flask, HTML</p>		
프로젝트 성과	<p>성공적으로 데모프로그램 구현. 이후 수업에서 부족한 데이터셋을 이후 새롭게 추가하여 소설 형태의 문장을 출력하도록 개선하였다.</p> <p>Github: https://github.com/dodoyeon/KoGPT2</p> <p>Slide: https://drive.google.com/file/d/1U6SUCUE42DXwjmQOJv1VSNaj4kZBVkah/view?usp=sharing</p>		

- **Github:** github.com/dodoyeon
- **Blog:** mari970.tistory.com