

Seminar Feedback Report (Week13)

2021312738 소프트웨어학과 김서환

13주차 세미나는 최원제 연구생의 발표로 시작되었으며, 주제는 대규모 언어 모델(LLM)의 추론 능력을 소형 언어 모델(SLM)으로 증류하여 embodied AI 태스크를 효율적으로 수행할 수 있도록 하는 DEDAR 프레임워크였다. 본 연구는 ICML에서 발표되었으며, 발표자는 기계 학습 연구실의 박사과정 학생으로, 연구실은 고급 강화 학습 기법과 다양한 응용 분야에 주력하고 있다. 연구의 핵심 내용은 MDP 기반의 체인 오브 사고(COT) 프롬프트와 2단계 정책 계층 구조를 결합하여 LLM의 지식을 SLM으로 효과적으로 이전하는 방법론을 제안하는 것으로, 이를 통해 제한된 환경에서도 강력한 성능을 발휘할 수 있음을 입증하였다고 설명하였다. 특히 AI Tutorial과 HabitatSim 같은 시뮬레이터를 활용하여 Alfred 환경에서 312개의 전문가 데이터셋으로 평가를 수행하였으며, LLM 기반 방법 대비 성공률과 목표 조건 성공률에서 각각 21.6%, 12.3%의 성능 향상을 달성했다고 발표하였다. 또한 Diddle 정책이 작은 모델 크기에도 불구하고 강력한 성능을 보였으며, 향후 도메인 이동 설정에서의 몇 샷 최적화 연구로 확장할 계획임을 강조해 주셨다. Embodied AI 분야에서 LLM의 추론 능력을 효율적으로 증류하는 새로운 방법론을 제시하였다는 점에서 연구의 의의가 크다고 생각되었다.

다음 발표는 배수영 연구생의 발표가 이어졌으며 NAACL 2025에 게재될 예정인 "A Context-Adapted Prompt Generation for Devising Zero-Shot Question Answering in Large Language Models" 논문을 주제로 정보 및 지능 시스템에 대한 발표를 진행하였다. 해당 연구실은 이지형 교수님의 지도 아래 컴퓨터 비전, 자연어 처리, 멀티모달, 경량 AI 등 다양한 연구를 수행하고 있다. 발표자는 LLM이 사회적 편향(bias) 문제와 질문 응답 작업에서의 신뢰성 한계를 보이는 현상에 주목하여, 질문의 모호성을 탐지하고 그에 적합한 프롬프트를 적용하여 편향된 답변을 줄이고 정확도를 높이는 DeCAP 방법론을 제안하였다고 발표하였다. 특히 DeCAP은 질문 모호성 감지 단계와 중립 답변 안내 생성 단계로 구성되며, BBQ와 Uncover 데이터셋을 사용한 평가에서 정확도와 편향 점수 측면에서 기존 방법들보다 뛰어난 성능을 보였다고 하셨다. BBQ 데이터셋에서는 정확도가 약 22.87%, 편향은 약 12% 개선되었으며, Uncover 데이터셋에서는 정확도가 약 50%, 편향은 약 2.27% 개선되었다고 발표하셨다. DeCAP이 추가 학습 없이도 다양한 편향 범주에서 일관된 성능을 보였다는 점을 강조하면서, 향후 더 넓은 도메인에서의 적용 가능성과 성능 개선에 대한 후속 연구의 필요성도 언급해 주셨다. 사회적 편향을 완화하고 LLM의 정확도를 동시에 높이는 새로운 접근법을 배울 수 있었던 유익한 세미나였다.