

Gen AI / LLM 활용 및 성능 개선 (AI Advanced)

1. 교육 일시

7/21(월) - 7/23(수) 전일 교육

2. 교육 목표 및 내용

교육 목표	<ul style="list-style-type: none">서비스에 필요한 Gen AI / LLM 관련 기술을 이해하고 적용할 수 있다.Gen AI 모델 학습 및 고도화를 통해 원하는 품질을 찾을 수 있다.
기대 효과	<ul style="list-style-type: none">Gen AI / LLM 을 서비스 관점의 적용 방안을 파악하고, 실서비스에 대한 개선이나 새로운 서비스를 개발할 수 있다.

3. 교육 진행

※ AM 10:00 - PM 7:00 까지 기본 교육 시간으로 진행됩니다.

• 오프라인 3일 과정

교육 시간표

이론+실습	프로젝트
-------	------

	Day 1	Day 2	Day 3
10:00 - 11:00	<ul style="list-style-type: none">- 개념/모델 이해<ul style="list-style-type: none">└ Deep learning└ Language model└ Transformer└ GPT1, GPT2, GPT3, ChatGPT (RLHF)└ Close Model / Open Model / Reasoning Model└ Model Evaluation	<ul style="list-style-type: none">- RAG<ul style="list-style-type: none">└ Loader└ Chunking└ Embedding Models└ Retrieval└ Vector Database└ Evaluation- PAC 성능 개선 / 최적화	<ul style="list-style-type: none">- Function Call 심화<ul style="list-style-type: none">└ Function Calling 기본 개념 및 이해└ Parallel Multiple Function Calls└ LLM Pipeline

11:00 - 12:00	<ul style="list-style-type: none"> Memory bound Optimization - On device AI <ul style="list-style-type: none"> Ollama sLLM Constraints (HW, Cost) Optimization: MoE, Kernel, Quantization 	- RAG 성능 개선 / 최적화 <ul style="list-style-type: none"> Query translation Decomposition Routing Query structuring Indexing Retrieval Graph RAG Contextual Retrieval 	<ul style="list-style-type: none"> AI Agent (LangGraph 활용) Integration - LLM 관련 오픈소스 소개
12:00 - 13:00			
13:00 - 14:00	- Prompt Engineering 주요 기법들 <ul style="list-style-type: none"> few shot, CoT, Reasoning Tokens, ReAct 	- 출처 내용 기반 답변 요약 기능	
14:00 - 15:00	- Prompt Engineering 심화 <ul style="list-style-type: none"> 타 서비스 시스템 프롬프트 분석 Prompt Format, Prompt Generator 프롬프트 비용, 프롬프트 압축 프롬프트 평가 	<Project : Perplexity Clone 서비스 개발> (RAG) <ul style="list-style-type: none"> 질의에 대한 출처 제공 출처 내용을 요약하여 답변 제공 	<Project : Perplexity Clone 서비스 개발> (AI Agent, 콘텐츠 요약 구현) <ul style="list-style-type: none"> LangGraph: GraphRAG / Agentic RAG 질문에 대한 출처를 제공 출처 내용 기반 답변 요약 기능 추가 (검색 추천 기능 구현) <ul style="list-style-type: none"> 출처 내용 기반 답변 요약 기능 추가 (프로젝트 결과 공유 및 선정 후 발표)
15:00 - 16:00			
16:00 - 17:00	(환경구축, Prompt Engineering) <ul style="list-style-type: none"> 프로젝트 설명 및 환경 구축 Model 연동 및 프로젝트 환경 구축 Prompt Engineering 적용 및 평가 Automated Prompt Generation 	<ul style="list-style-type: none"> Finetuning 원리/기법 이론 설명 Finetuning vs RAG finetune 등 관련 기술들 (LoRA, PPO, trl) OpenAI Fine-tuning API 경량 모델 파인튜닝 (1h) - 평가 <ul style="list-style-type: none"> 정량 평가 (Rouge, BLEU, Benchmarks Datasets) 정성 평가 (human eval, metrics) 	
17:00 - 18:00			- Multimodal 이해 <ul style="list-style-type: none"> 적용된 서비스/사례 소개/학습법 OpenAI API 연동 서비스
18:00 - 19:00			- AI 서비스 개발 방향성 (LLMOps stack) <ul style="list-style-type: none"> Q&A Closing

※ Perplexity AI : LLM을 사용하여 질문에 대한 정확한 답변을 제공하는 AI 기반 검색 엔진