

# 임베디드 환경 최적화 LLM 개발 스터디 프로젝트 로드맵

## 1. 프로젝트 개요 및 목적

임베디드 환경에 최적화된 로컬 LLM(Local Large Language Model) 개발을 목표로 한 스터디 프로젝트입니다. 대형 언어 모델은 주로 클라우드 GPU 환경에서 동작하지만, 개인정보 보호, 오프라인 사용, 지연 시간 단축 등을 위해 엣지 디바이스(임베디드 기기)에서 구동하려는 수요가 높아지고 있습니다 ①. 본 프로젝트는 중급 수준의 머신러닝 스터디 멤버들이 이러한 요구에 부응하는 기술을 학습하고, 직접 경량화된 LLM을 구현하여 임베디드 디바이스에서 구동해 보는 것을 목적으로 합니다. 이를 통해 최신 NLP/딥러닝 트렌드인 대규모 언어 모델을 자원 제약 환경에서 다루는 경험을 쌓고, 모델 최적화 기술과 엣지 AI 개발 역량을 향상시키는 것이 궁극적인 목표입니다.

### 왜 임베디드 환경에서 LLM인가?

LLM을 로컬 디바이스에서 실행하면 인터넷 연결 없이도 동작하여 네트워크 지연을 줄이고, 민감한 데이터를 외부로 보내지 않아 프라이버시를 강화할 수 있습니다. 또한 클라우드 비용 절감과 규제 준수 측면에서도 이점이 있어, 최근 많은 연구가 경량화된 LLM의 엣지 실행에 집중되고 있습니다 ①.

## 2. 전체 개발 로드맵 개요

프로젝트는 5단계 로드맵으로 진행됩니다. 각 단계마다 학습 및 개발 활동과 산출물이 정의되어 있으며, 단계가 거듭될수록 점진적으로 실제 임베디드 환경에서 동작하는 LLM 시스템을 완성하게 됩니다. 전체 흐름은 다음과 같습니다:

1. 리서치 단계 - LLM 경량화와 임베디드 AI 관련 연구 조사 및 기획 수립
2. Kaggle 실습 단계 - Kaggle 경진대회 및 노트북을 활용한 실전 모델 개발 경험 축적
3. 모델 경량화 및 튜닝 단계 - 선택한 LLM의 경량화(압축) 및 파인튜닝 진행
4. 엣지 디바이스 포팅 단계 - 경량화된 모델을 임베디드 하드웨어에 배포 및 구동
5. 최적화 및 멀티 플랫폼 확장 단계 - 모델의 성능 최적화와 다양한 플랫폼으로의 확대 적용

각 단계별로 구체적인 활동 내용, 사용 도구, 산출물, 그리고 성공 기준(완료 정의)을 아래에 상세히 설명합니다. 이 로드맵을 따라가면, 초기 아이디어 구상부터 최종 임베디드 장치 상용 수준 구현까지 체계적으로 접근할 수 있습니다.

## 3. 단계별 세부 설명과 산출물

### 3.1 1단계 - 리서치 (Research)

**내용:** 프로젝트의 기초를 다지는 단계로, LLM 모델 경량화 기법과 임베디드 AI 사례에 대한 문헌 연구를 수행합니다. 최신 학술 논문, 기술 블로그, 오픈소스 프로젝트를 조사하여, 자원 제약 환경에서 LLM을 구현하기 위한 핵심 기술 요소를 파악합니다. 예컨대, 모델 양자화(quantization), 프루닝(pruning), 지식 증류(distillation), 메모리 최적화 기법 등이 어떻게 LLM에 적용되는지 조사합니다 ②. 또한 Hugging Face나 PyTorch 허브를 통해 경량 모델 사례(예: DistilBERT, TinyLlama 등)를 실험적으로 분석하고, 프로젝트에 적합한 베이스 모델 선택 기준도 수립합니다.

- **사용 도구:** ArXiv 논문, 주요 학회 튜토리얼, 관련 서베이 논문, 기술 블로그, Hugging Face 모델 저장소, etc. (예: "A Comprehensive Survey of Small Language Models..." 논문이나 Analytics Vidhya 블로그 등 최신 자료)
- **산출물:** 조사 내용 정리 보고서 또는 위키 페이지. 여기에는 유망한 경량화 기법 목록, 후보 베이스 LLM 모델 비교표, 활용 가능한 오픈소스 툴킷 정보 등이 포함됩니다. 또한 단계 3부터 적용할 모델 및 기법 선택안과 전체 개발 전략 초안도 산출합니다.

- **성공 기준:** LLM 경량화와 임베디드 적용에 대한 **이해도 향상** 및 팀 내 **공유**. 구체적으로, 최소 **3~5건 이상의 핵심 기술**을 파악하고 정리했으며, 프로젝트에서 사용할 **베이스 LLM 모델**을 (예: LLaMA 2 7B, GPT-J 계열 등) 선정 완료한 상태를 목표로 합니다.

### 3.2 2단계 – Kaggle 실습 (Kaggle Practice)

**내용:** 리서치 단계에서 얻은 이론과 아이디어를 실제로 적용해보기 위해 **Kaggle 플랫폼에서 실습**을 진행합니다. Kaggle에는 LLM 및 NLP 관련 다양한 **경진대회**와 **공개 노트북**이 존재하므로, 이를 활용하여 팀원들이 **모델 개발 경험**을 쌓습니다. 예를 들어, LLM을 활용한 분류나 질의응답 대회를 선택해 참가하거나, Kaggle 노트북 환경에서 작은 언어 모델을 **Fine-tuning**하고 **추론**까지 수행해 봅니다. Kaggle의 제한된 자원(예: GPU 메모리 제약) 하에서 모델을 다루는 경험은 이후 임베디드 환경의 제약 대응에 큰 도움이 됩니다 <sup>3</sup>. 또한 토론 포럼을 통해 **전 세계 참가자들의 아이디어**를 접하고, **데이터 전처리, 모델 앙상블, 성능 평가** 등 실무 팁을 배울 수 있습니다.

- **사용 도구:** Kaggle Notebooks (무료 GPU), Python ML 스택(Pandas, PyTorch/TensorFlow, Hugging Face Transformers), Kaggle 대회 데이터셋. 필요시 Kaggle Discussions 및 외부 커널 참고.
- **산출물:** 선택한 경진대회에서의 **노트북 코드, 모델 학습 결과물(예: 제출 파일)**, 그리고 실습 과정을 통한 **교훈 정리** 문서. 예컨대 “어떤 기법이 효과적이었고 무엇이 어려웠는지”에 대한 팀 회고 자료를 만듭니다.
- **성공 기준:** 팀원 전원이 Kaggle 실습을 통해 **실제 모델 학습/튜닝/평가 사이클**을 경험하는 것이 목표입니다. 가능하다면 **Kaggle 대회 순위 상위권 진입**이나 **논의 게시판에 아이디어 공유**를 해보는 것도 좋습니다. 최소한 베이스라인 모델을 개선하여 **스코어 향상**을 이루고, 실습 전후로 팀원들의 **실전 역량 향상**이 확인되면 성공입니다.

### 3.3 3단계 – 모델 경량화 및 튜닝 (Model Lightweighting & Fine-tuning)

**내용:** 본격적으로 **선정한 LLM 모델을 경량화**하고 과제에 맞게 **튜닝**하는 단계입니다. 우선 리서치한 기법 중 프로젝트에 적합한 것을 선택하여 적용합니다. 일반적인 접근은 **사전 학습된 대형 모델을 최대한 작은 크기로 압축**한 후, **파인튜닝**을 통해 성능을 보완하는 것입니다. 구체적인 기법들은 다음과 같습니다:

- **양자화(Quantization):** 모델 가중치의 표현 정밀도를 낮추어 메모리 Footprint를 줄입니다. 예를 들어 32-bit 부동소수점을 8-bit 정수로 양자화하면 **메모리 사용을 4배 축소**할 수 있고, 4-bit로 하면 **8배 축소**가 가능합니다 <sup>4</sup>. 실제 연구에서도 8억 파라미터 규모 LLM을 16-bit에서 4-bit로 양자화하여 모델 크기를 약 16GB→4GB로 줄이고도 성능 저하를 최소화한 사례가 있습니다 <sup>5</sup>. 양자화는 **추론 속도 향상**에도 도움이 되며, Kaggle 등 제한된 환경이나 임베디드 디바이스에서 **LLM을 구동 가능하게 만드는 핵심 기술**입니다. 양자화 적용 시에는 Hugging Face의 `bitsandbytes` (int8 지원)나 GPTQ, NVIDIA TensorRT-LLM 등의 툴을 활용할 수 있습니다. 또한 **혼합 정밀도나 레이어별 다른 비트수 사용** 등 고급 기법도 고려합니다 <sup>4 6</sup>.
- **지식 증류(Knowledge Distillation):** 매우 큰 **교사(teacher) 모델**의 지식을 **경량 학생(student) 모델**에 증류하는 방법입니다. 교사 LLM의 출력(예측 분포)을 정답처럼 활용하여 작은 모델을 학습시키면, 학생 모델이 교사의 성능 일부를 학습하여 훨씬 **파라미터 수가 적어도 높은 성능**을 내도록 할 수 있습니다. 예를 들어 MiniLLM 등의 연구에서는 거대 언어모델들을 수백만~수억 파라미터 대 모델로 증류하여 **긴 텍스트 생성 품질을 유지**하는 데 성공했습니다 <sup>7</sup>. 증류를 위해서는 학습 데이터와 교사 모델이 필요하며, 자체 태스크 데이터 외에도 **생성 데이터 추가**를 고려할 수 있습니다.
- **모델 구조 프루닝(Pruning) 및 구조 변경:** 모델의 일부 **가중치나 뉴런을 제거**하여 경량화하는 기법입니다. 비중요한 연결을 0으로 만들어 **희소화**하거나, **어텐션 헤드나 레이어**를 제거하는 등의 **구조적 프루닝**을 적용합니다 <sup>8</sup>. 다만 LLM에 선불리 프루닝을 적용하면 성능 저하가 크므로, 필요 시 **미세조정(fine-tuning)**과 병행하여 영향 최소화를 도모합니다. 또는 기반 모델로 애초에 경량 구조(BERT→DistilBERT처럼)를 채택하는 것도 방법입니다.
- **효율적 파인튜닝 기법:** 경량화와 동시에 원하는 태스크에 맞게 모델을 최적화해야 합니다. 파인튜닝 시 모든 파라미터를 다 학습시키기보다 **LoRA**와 같은 **경량 퍼램추가(PEFT)** 기법을 활용하면 메모리와 연산량을 아낄 수

있습니다. LoRA는 대형 모델의 **일부 행렬에 저랭크 행렬**을 추가학습하는 것으로, 예를 들어 Kaggle LLM Science Exam 솔루션에서도 7B~13B 모델의 모든 **선형계층만 LoRA로 미세조정**하여 성능을 향상시켰습니다 <sup>9</sup>. 또한 최근 각광받는 **QLoRA** 기법은 4-bit 양자화를 병행한 LoRA 미세조정으로, 적은 GPU 메모리로도 거대 모델(예: 70B)을 튜닝 가능하게 해줍니다 <sup>10</sup>. 이러한 기법들을 활용하면 임베디드 대응을 위해 **모델 용량을 줄이면서도 필요한 성능을 확보**할 수 있습니다.

- **평가 및 반복:** 경량화/튜닝 결과 모델의 성능을 **검증 데이터셋**으로 평가합니다. 만약 성능 저하가 너무 크면, 양자화 비트수를 높이거나, 프루닝 정도를 완화하거나, 추가 튜닝(Epoch 증가 등)을 통해 **트레이드오프 조절**을 반복합니다. 이때 평가 지표는 프로젝트 목적에 맞게 선택하며 (예: 답변 정확도, 생성 품질 점수, 속도(FPS) 등), 임베디드 최종 적용을 고려해 **메모리 사용량, 추론 속도**도 함께 측정합니다.
- **사용 도구:** PyTorch 또는 TensorFlow 딥러닝 프레임워크, Hugging Face Transformers + Accelerate, PEFT 라이브러리, BitsandBytes (INT8), GPTQ 등 양자화 킷, Distiller (증류용 툴) 등. 또 실험 추적을 위해 Weights & Biases 같은 툴을 사용할 수 있습니다.
- **산출물:** 경량화된 최종 **모델 가중치** 파일 (예: FP16 대비 메모리 X배 감소된 모델), 해당 모델의 **성능 평가 보고서 및 비교표**. 또한 각 기법 적용에 따른 효과를 정리한 **실험 노트**도 산출합니다. 예를 들어 "INT8 양자화 적용 시 정확도 1% 하락, 속도 1.5배 향상" 등의 결과를 기록합니다. 필요하다면 추후 배포를 고려해 **모델 카드 (Model Card)** 형태로 모델 정보를 정리합니다.
- **성공 기준:** 원래 모델 대비 **상당한 경량화**(예: 모델 크기 25% 이하로 감소) 달성 및 **목표 성능 수치 유지**가 기준입니다. 구체적으로, 선택한 평가 지표에서 **baseline 대비 ~90% 이상의 성능을 유지**하면서도 메모리 사용과 추론 속도가 크게 개선된 모델을 얻는 것을 목표로 합니다. 예를 들어 10GB 모델을 2GB로 줄이고도 핵심 과제 정확도를 95% 이상 달성했다면 성공으로 간주합니다. 또한 이 단계까지 완료하면 다음 임베디드 배포를 위한 **준비된 모델**이 확보됩니다.

### 3.4 4단계 – 엣지 디바이스 포팅 (Edge Device Porting)

**내용:** 경량화된 모델을 실제 **임베디드 하드웨어에 배포**하여 동작시키는 단계입니다. 여기서 임베디드 디바이스란 **라즈베리 파이, NVIDIA Jetson, 안드로이드 폰, 임베디드 보드** 등 제한된 자원을 가진 모든 로컬 장치를 포함합니다. 우선 현재 모델을 해당 디바이스에서 효율적으로 구동할 수 있는 **형태와 런타임**으로 변환합니다. 주요 포팅 전략은 아래와 같습니다:

- **모델 포맷 변환:** 딥러닝 프레임워크에서 학습한 모델을 범용적이고 최적화된 포맷으로 변환합니다. 대표적으로 **ONNX**로 변환하면 다양한 플랫폼에서 가속 실행이 가능합니다. ONNX Runtime은 CPU, GPU, NPU 등 **다양한 하드웨어 가속자**를 지원하고 불필요한 연산을 제거하는 **그래프 최적화**를 제공하여 엣지 추론을 최적화합니다 <sup>11</sup>. 또는 **TensorFlow Lite**로 변환하면 Android나 임베디드 Linux 환경에서 **경량 런타임**으로 구동할 수 있고, 선택적 연산자만 포함하여 바이너리 크기를 줄이거나 Edge TPU delegation 등의 **디바이스 특화 최적화**를 적용할 수 있습니다 <sup>12</sup>. 변환 과정에서 INT8 양자화 최종 적용(PTQ)도 함께 수행하여 모델 파일 크기를 줄입니다.
- **플랫폼별 실행 환경 설정:** 대상 디바이스에 맞는 라이브러리 세팅을 합니다. 예를 들어 Raspberry Pi같이 CPU만 있는 경우 **임베디드 BLAS** 라이브러리 최적화나, `llama.cpp`처럼 **C++ 기반 경량 추론엔진** 사용을 고려합니다 <sup>13</sup>. Llama.cpp는 LLM을 위해 특화된 C++ 구현으로, 메모리 맵핑을 통해 RAM 사용을 줄이고 CPU SIMD 최적화로 수십억 파라미터 모델도 노트북이나 스마트폰에서 동작시킨 유명 프로젝트입니다 <sup>13</sup>. 이러한 구현은 파이썬보다 메모리 및 속도 면에서 유리하므로, 필요한 경우 모델 가중치를 변환(예: ggml 포맷)하여 사용합니다. Android의 경우 Android NNAPI나 GPU/OpenCL 가속을 활용할 수 있고, iOS라면 CoreML 변환을 고려합니다.
- **디바이스 리소스 튜닝:** 모델을 디바이스에서 로드하고 **실제 inference를 테스트**합니다. 이때 디바이스의 **메모리 한계, 저장공간, 열 발생, 전력 소비** 등을 모니터링하여 문제가 없는지 확인합니다 <sup>14</sup> <sup>15</sup>. 시퀀스 길이가 긴 LLM의 경우 **KV 캐시(Key-Value Cache)** 메모리가 크게 작용하므로, 필요하면 **KV 캐시까지 양자화**하거나

순차 토큰 생성을 **캐싱하여 속도 개선**하는 등 추가 최적화를 적용합니다 <sup>16</sup> <sup>17</sup> . 예를 들어, KV 캐시를 활용하면 매 토큰 생성 시 이전 토큰에 대한 **중복 계산을 피하여** 추론 복잡도를 크게 줄일 수 있어, 긴 문장 생성 시 **속도를 선형으로 개선**합니다 <sup>18</sup> .

- **사용 도구:** ONNX 변환 (`torch.onnx` 또는 `tf2onnx` ), ONNX Runtime C++/Python, TensorFlow Lite converter 및 Java/C++ API, llama.cpp 빌드 스크립트, Embedded Linux 툴체인, Android Studio (JNI 연결 시), etc. 또한 하드웨어 모니터링 툴 (tegrastats, Android Profiler 등)을 사용합니다.
- **산출물:** 임베디드 디바이스에서 동작하는 **프로토타입 애플리케이션**. 예를 들어 Raspberry Pi에서 터미널 기반 Q&A 챗봇이 돌거나, Android 폰에서 간단한 데모 앱으로 LLM 답변을 생성하는 모습. 이와 함께 **배포 가이드 문서**를 작성하여, 해당 디바이스에서 필요한 세팅과 실행 방법, 성능 측정 결과를 정리합니다 (예: "Jetson Nano에서 ONNX Runtime으로 BERT QA 모델 실행 - 메모리 1GB 사용, latency 2초" 등의 정보).
- **성공 기준:** 선택한 **타겟 임베디드 환경에서 모델이 원활히 구동**되는 것이 핵심 성공 기준입니다. 구체적으로, (i) 모델이 메모리 초과 등 없이 로드되고, (ii) 한 번 추론(예: 답변 생성 또는 분류)이 **합리적인 시간 내(수 초 내)** 완료되며, (iii) 결과가 유의미하게 나온다면 성공입니다. 추가로 디바이스의 열이나 응답성 등 **실사용 가능성**도 평가하여 큰 문제가 없을 것이라 판단되면 해당 단계 완료로 간주합니다.

### 3.5 5단계 - 최적화 및 멀티 플랫폼 확장 (Optimization & Multi-platform Expansion)

**내용:** 마지막으로, 구현한 LLM 시스템을 **성능 최적화**하고 **여러 플랫폼으로 확장**하는 단계입니다. 4단계에서 단일 환경 구동에 성공했다면, 보다 **실전 배포 수준의 개선**을 통해 활용도를 높입니다.

- **성능 프로파일링 및 최적화:** 임베디드 디바이스에서 프로파일링을 실시하여 **병목 요소**를 파악합니다. 예를 들어 CPU 사용률, 메모리 대역폭, 특정 연산(Fused operations 가능 여부) 등을 분석합니다. 이 정보를 바탕으로 스레드 병렬화 설정 조정, XLA 컴파일러 최적화, 메모리 풀링 등 **세부 최적화**를 적용합니다. 또한 가능하다면 **하드웨어 가속기** 사용을 극대화합니다 (예: DSP/NPU가 있다면 해당 연산 위임). GPU가 있는 엣지 장치의 경우 TensorRT 등의 **벤더 최적화 runtime**을 적용해 볼 수 있습니다.
- **멀티 플랫폼 지원:** 초기 대상 외에 **다른 플랫폼으로의 이전**을 시도합니다. 예를 들어 원래 Linux 기반 보드에 배포했다면 **모바일(Android/iOS)** 버전으로도 포팅해보고, x86 데스크톱에서도 **동일 모델이 동작**하게 하여 유지 보수합니다. ONNX 같은 공통 포맷을 사용했다면 비교적 수월하게 여러 환경에서 구동할 수 있습니다. 각 플랫폼별로 상이한 제약 (모바일은 배터리/발열, PC는 쓰레딩 최적화 등)이 있으므로 이에 맞게 설정을 변경합니다. 또한 멀티 플랫폼 테스트를 통해 코드베이스를 정리하고 **CI/CD** 파이프라인으로 다양한 환경에 배포하는 자동화도 고려할 수 있습니다.
- **기능 및 UX 개선:** 모델의 응답 속도가 최적화되었는지 최종 확인하고, 필요시 **속도를 높이기 위해 추가 양자화** (예: 4-bit->3-bit 시도)나 **지연 로드** 등의 테크닉을 씁니다. 동시에 사용자 입장에서 쓸모있게 **응용 수준 기능**을 개선합니다. 예를 들어 질의응답 챗봇이라면 대화 기록 유지, 간단한 GUI, 또는 음성 입출력 연동 등을 붙여볼 수 있습니다. **모델 정확도 개선**을 위해 추가 파인튜닝이나 데이터 보강을 할 수도 있습니다. 이 단계에서는 팀의 관심사에 따라 프로젝트 산출물을 **제품 수준으로 가다듬는 작업**을 진행합니다.
- **사용 도구:** 프로파일러(PyTorch Profiler, Android systrace 등), 다양한 플랫폼 빌드 도구(Docker 멀티스테이지 빌드, cross-compile 툴 등), CI 도구(GitHub Actions, Jenkins) 등. 성능 최적화에는 low-level 코드 작성(C++), SIMD 인스트럭션 활용도 포함될 수 있습니다.
- **산출물:** **최적화 보고서** (프로파일링 결과와 적용한 최적화 리스트, 그에 따른 성능 향상 수치), 멀티 플랫폼 동작 **스크린샷 또는 동영상** (예: Android 앱과 Raspberry Pi 양쪽에서 구동되는 모습). 그리고 사용자가 참고할 **튜토리얼 문서나 README**도 작성하여, 전체 시스템을 재현하거나 설치할 수 있는 가이드로 제공합니다.
- **성공 기준:** 모델이 초기 대비 **두드러지게 최적화**되고, **여러 환경에서 일관되게 동작**하면 성공입니다. 예를 들어 원래 추론에 5초 걸리던 것을 최적화로 2초로 단축, 메모리는 500MB 줄이고, 추가로 Android 폰에서도 구동

확인 — 이런 식의 결과라면 최종 목표 달성입니다. 또한 팀원들이 프로젝트를 통해 임베디드 LLM 개발에 자신감을 얻고, 산출물을 공개하거나 내부 공유하여 **프로젝트를 마무리**하면 성공적입니다.

## 4. Kaggle 주요 경진대회 사례 요약 (SLM/LLM 관련)

스터디 중 참고할 만한 **최근 Kaggle 경진대회**들을 몇 가지 소개합니다. 이들은 **Small/Large Language Model**을 다루는 흥미로운 문제들을 제공하여, 우리 프로젝트와 관련된 통찰을 얻거나 실습으로 활용하기 좋습니다.

- **LLM Science Exam (2023)** – 목적: GPT-3.5 등 **대형 언어모델이 만든 과학 시험 문제**(지문 기반 객관식)에 대한 **정답 추론** <sup>19</sup>. 참가자들은 5지선다 문제에 대해 각 선택지가 정답일 확률을 예측해 순위를 매기는 과제를 풀었습니다. 활용 포인트: **Retrieval-Augmented Generation (RAG)** 기법과 **작은 LLM 활용**이 크게 주목받았습니다. Kaggle 노트북 환경 제약으로 약 10B급 모델까지 사용 가능했기에, **7B~13B 수준의 비교적 작은 LLM들을 앙상블**하고 위키피디아 검색 등 **외부 지식검색**을 결합하여 거대 모델이 만든 문제를 푸는 접근이 유효했습니다 <sup>3</sup>. 상위권 솔루션들을 통해 **LoRA를 이용한 부분 파인튜닝** <sup>9</sup>, **QLoRA로 4-bit 양자화 파인튜닝** <sup>20</sup>, **다양한 재순회화 기법** 등 우리 프로젝트에 응용 가능한 아이디어를 얻을 수 있습니다.
- **LLM – Detect AI Generated Text (2023)** – 목적: 주어진 에세이가 **학생(인간)이 작성한 것인지 LLM이 생성한 것인지**를 분류하는 모델 개발 <sup>21</sup>. 실제 교육 현장의 글과 AI 산출물을 섞어놓고 이를 구분하는 문제가 출제되었습니다. 활용 포인트: **텍스트 분류를 위한 경량 언어모델 활용**과 **데이터 불균형 문제 대응**이 주요 과제였습니다. 훈련 데이터에서 인간 작성 글이 압도적으로 많고 AI 작성 글은 적었기 때문에, 참가자들은 ChatGPT 등으로 **추가 AI 생성 데이터를 보강**하거나, **특징 공학+전통 ML vs 파인튜닝 언어모델**을 비교하며 접근했습니다 <sup>22</sup> <sup>23</sup>. 상위 솔루션에서는 DeBERTa 등 **사전학습 언어모델을 미세조정**하여 높은 정확도를 거두었으며, 이는 우리 프로젝트에서 **소형 LLM의 분류 정확도**를 끌어올리는 팁이 될 수 있습니다.
- **H2O.ai Predict the LLM (2023)** – 목적: 주어진 텍스트 응답을 생성한 **LLM의 정체**를 맞추는 것입니다. 7개의 서로 다른 LLM이 동일 질문에 답한 응답 데이터가 주어지면, 각 응답이 어떤 모델(GPT-3, Falcon, LLaMA 등)에서 나온 것인지 맞추는 **모델 식별** 과제였습니다. 활용 포인트: **각 LLM의 출력 스타일 차이**를 학습하는 흥미로운 문제로, **텍스트의 통계적 특징, 어휘 사용, 문체** 등을 분석해야 했습니다. 이 대회 목표 자체가 “**7개 중 어떤 LLM의 산출물인지 검출**”이라고 명시될 정도로 특이하며 <sup>24</sup>, 이를 통해 **모델 간 특성 차이**, 나아가 **LLM 수탁(출처) 검증 기법**을 엿볼 수 있습니다. 상위권 풀이를 보면 클래식한 **자연어 특징공학+ML** 접근과 **Transformer 파인튜닝** 접근이 혼합되었고, Hugging Face의 **H2O LLM Studio** 등을 활용해 모델을 Fine-tune하여 성능을 높인 사례가 있었습니다. 우리 팀이 만약 여러 다른 모델을 비교하거나, 한 모델의 출력을 다른 경량 모델로 모방(종류)할 때 시사점을 줄 수 있는 대회입니다.
- **LLM Prompt Recovery (2024)** – 목적: 주어진 **출력 텍스트에 대해, 원래 LLM에 입력된 프롬프트를 역추론**하는 도전적인 과제 <sup>25</sup>. 예컨대 어떤 변환(스타일 변화나 요약 등)이 적용된 문장이 주어지면, 그 변환을 일으킨 원본 프롬프트(명령)를 추측해야 합니다. 활용 포인트: 프롬프트 엔지니어링과 **역방향 추론**을 다루는 실험적인 문제로, LLM의 동작을 **이해/해석하는 능력**이 필요합니다. 2024년 Kaggle에서 상당히 큰 상금(\$200k)과 함께 열려 주목받았으며, 다양한 **NLP 알고리즘, 언어 모델, 취약점 탐지 기법**이 동원되었습니다. 이를 통해 프롬프트와 출력 간 매핑을 학습하는 접근(예: **문장 임베딩 및 코사인 유사도 활용**) 등을 엿볼 수 있고, LLM의 내부 표현에 대한 **이해도**를 높이는 연습이 됩니다. 우리 프로젝트와 직접적 연관은 적지만, **LLM 활용 범위 확장 및 보안** 관점에서 참고될 수 있습니다.
- **LLMs – You Can’t Please Them All (2024)** – 목적: LLM을 **에세이 심사자(judge)**로 사용할 때, **악의적인 입력(advserial input)**으로 그 평가를 속일 수 있는지 찾아내는 대회입니다. 즉, 주어진 자동 채점 LLM 시스템에 최대한 높은 점수를 받도록 속이는 에세이를 생성하거나, LLM 평가의 취약점을 드러내는 것이 목표였습니다. 활용 포인트: **LLM의 안전성 및 평가 안정성**과 관련된 흥미로운 챌린지로, 대회 참가자들은 LLM 판정 알고리즘이 오작동하도록 하는 **프롬프트 기법, 내용 교란** 등을 탐구했습니다. 이를 통해 LLM 기반 시스템의 **공격과 방어 전략**을 배울 수 있고, 우리 팀에는 LLM을 응용할 때 고려해야 할 **윤리적/안전성 측면**에 대한 토론 소재가

되었습니다. (※ 이 대회는 결과보다도 문제를 통해 "LLM의 판단 신뢰도"라는 주제를 환기시켰다는 점에서 의의가 있습니다.)

以上的 경진대회들은 **실제 문제를 통한 LLM 활용**을 다양하게 보여주므로, 스터디 진행 중 해당 대회 페이지의 **문제 설명, 데이터, 논의글**을 읽어보고 아이디어를 얻거나, 직접 코드 실행을 따라해 보면 좋습니다. 특히 상위 입상 솔루션의 요약글을 읽어보면 **경량 모델로도 창의적인 접근**을 한 예시가 많아, 프로젝트에 인사이트를 줄 것입니다.

## 5. 팀 운영 팁과 기술 전략

### 팀 운영 팁 (협업 및 학습 전략)

- **역할 분담과 정기 공유:** 프로젝트 로드맵 단계별로 팀원들에게 역할을 분담하여 **병렬로 진행**합니다. 예를 들어 리서치 담당, Kaggle 실습 담당, 모델 튜닝 담당, 디바이스 테스트 담당 등을 정해 각자 deep-dive한 뒤 **위클리 미팅**에서 지식을 공유합니다. 이렇게 하면 한정된 기간 내에 폭넓은 분야를 다룰 수 있고, 모두가 모든 부분을 배우는 효과도 얻을 수 있습니다.
- **애자일 방식 진행:** 짧은 주기의 **스프린트**를 구성하여 매주/격주 목표를 설정하고 회고합니다. 각 단계별 목표를 **마이크로 태스크**로 쪼개고 진행 상황을 Kanban 보드 등에 시각화하면 추적이 수월합니다. 특히 3~5단계에서는 예기치 않은 기술 난관이 등장하기에, 유연하게 일정과 계획을 조정할 수 있도록 **애자일**한 태도가 중요합니다.
- **버전 관리 및 협업 도구 활용:** 소스코드와 실험 결과는 모두 **Git**으로 **버전관리**하고, 모델 가중치나 데이터도 필요하면 **Git LFS**나 Kaggle Datasets, 구글 드라이브 등을 활용해 관리합니다. 이때 **실험 노트(예: Jupyter Notebook)**도 일관되게 관리하여, 누가 어떤 실험을 했고 결과가 어땠는지 모두가 볼 수 있게 합니다. 또한 Slack/Discord 등으로 상시 소통 채널을 운영하고, 중요한 의사결정이나 아이디어는 문서화하여 공유합니다.
- **학습 리소스 적극 공유:** 중간에 각자가 찾아낸 유용한 **튜토리얼, 강의, 깃헙 리포지토리** 등을 팀 내 채널에 공유하세요. 예를 들어 “양자화 설명이 잘 된 블로그”나 “ONNX 최적화 예제 코드” 등을 함께 스터디하면 개별 시행착오를 줄이고 팀 전체 역량을 높일 수 있습니다.
- **리스크 관리:** 프로젝트 진행 중 **기술적 난이도**나 **시간 부족** 리스크가 발생할 수 있습니다. 예를 들어 “생각보다 임베딩 보드에서 모델이 안 돌아간다”와 같은 문제가 생기면, 로드맵을 재검토하여 **대안 경로**를 마련합니다. (모델 크기를 더 줄인다든지, 다른 보드로 바꾼다든지 유연하게 대처) 애초에 리스크 요인을 목록화하고 사전에 대비책을 생각해두면 좋습니다.

### 기술 전략 (개발 시 고려할 사항)

- **모델 선택 전략:** 처음부터 모든 것을 직접 만들기보다, **검증된 사전학습 LLM**을 기반으로 시작하는 것이 현명합니다. 현재 공개된 LLM 중 소형 모델로 성능이 우수한 것들을 검토합니다. 예를 들어 **LLaMA-2 7B, Mistral 7B, GPT-J 6B, Flan-T5 (Small/XL)** 등이 후보가 될 수 있습니다. 선택 시에는 모델의 라이선스(Open/Open-source 여부), 한국어 지원 여부(필요하다면), 커뮤니티 지원 등을 고려합니다. 또한 동일 파라미터 규모라도 아키텍처에 따라 속도가 다를 수 있으므로(MHA 구조 최적화 등), **실제 벤치마크 결과**를 찾아보고 결정합니다.
- **데이터 및 태스크 설정:** 우리 프로젝트의 목적에 맞는 **구체적 태스크**를 정의하고 이에 맞는 데이터를 확보해야 합니다. 예를 들어 임베딩 기기에서 사용할 AI 비서라면 **멀티턴 대화 데이터**, 특정 산업 기기에 쓸 모델이라면 해당 도메인 코퍼스 등이 필요합니다. 공개 데이터셋 (예: KorQuAD, AI Hub 데이터 등)을 활용하거나 Kaggle 실습 단계에서 쓴 데이터를 재사용합니다. 경우에 따라서는 **데이터 증강**이나 **자체 크롤링**도 고려하여, 튜닝에 충분한 데이터를 마련합니다.
- **경량화 기술 조합:** 모델 경량화는 앞서 설명한 여러 기법을 **복합적으로 적용**하면 시너지가 큼니다 <sup>2</sup>. 양자화는 가장 손쉽고 효과적인 방법이며, 여기에 지식 증류를 더하면 성능저하 최소화에도 좋습니다. 필요하다면 프루닝으로 크기를 더 줄이고, 최종적으로 LoRA로 파인튜닝을 해서 성능 보충을 합니다. 이처럼 **양자화 + 증류 + LoRA** 등을 조합하면 개별 적용보다 성능 보존에 유리합니다. 실제 연구에서도 3비트까지 양자화하면서도 교사 모델로 보정하여 성능을 유지하는 등 다양한 시도가 있으므로, 팀 내에서도 아이디어를 내보고 실험적으로 검증해봅니다.

- **테스트 주도 개발:** 모델 개발 및 튜닝 과정에서 **단계별 테스트**를 구축하면 나중에 문제를 좁혀서 해결하기 좋습니다. 예를 들어 간단한 검증 질문 셋을 만들어서, 현재 모델이 얼마나 답을 잘하는지 또는 얼마나 빠른지를 지속적으로 측정합니다. 또 임베딩 포팅 전후에 동일 입력에 대한 출력을 비교하여 포팅 오류를 체크합니다. 이러한 테스트 자동화는 프로젝트 막바지에 **회귀(regression) 발생을 방지**하고 품질을 보장하는 데 도움이 됩니다.
- **포팅 및 최적화 전략:** 임베딩 장치별로 최적화 포인트가 다르므로, **플랫폼 특화 전략**을 세웁니다. 예를 들어 **ARM CPU** 기반 보드에서는 Neon SIMD를 활용하도록 컴파일 옵션을 주고, **GPU** 있는 장치는 배치 크기를 키워 병렬성을 높입니다. ONNX Runtime의 **Execution Provider**를 바꿔가며 (CUDA, TensorRT, OpenVINO 등) 테스트해보고 가장 빠른 것을 고릅니다. 메모리가 극히 적다면 **스트리밍 처리**나 **모델 쪼개기(모듈화)**도 고려합니다. 또한, llama.cpp 같이 **커뮤니티 최적화 프로젝트**를 적극 활용하면 많은 저수준 최적화를 직접 하지 않아도 되어 시간을 절약할 수 있습니다.
- **모델 평가 및 모니터링:** 최종 시스템에 대해 **다각적인 평가**를 수행합니다. 성능(정확도)과 속도 뿐 아니라, 임베딩 사용 시 중요한 **안정성(크래시 여부)**, **장시간 동작시 메모리 누수**, **응답 일관성** 등을 점검합니다. 필요하다면 텔레메트리 코드를 삽입해 메모리 사용량, 처리 시간 등을 로깅하고, 여러 번 반복 구동하여 변화를 관찰합니다. 이러한 모니터링을 통해 현장에서 생길 수 있는 문제를 선제 대응합니다.
- **윤리 및 사용자 피드백:** LLM을 실제 활용할 때는 **윤리적 이슈**도 간과할 수 없습니다. 임베딩 시라 하더라도 부적절한 응답을 내지 않도록 **필터링**이 필요하고, 사용자 데이터가 축적될 경우 **프라이버시 보호**에도 신경써야 합니다. 팀 내부적으로 이러한 이슈도 토의하고, 가능한 대책 (예: 반응형 프롬프트 필터, on-device 데이터 암호화 저장 등)을 마련합니다. 또한 최종 사용자의 피드백을 받아 개선점을 도출하는 **사용자 테스트** 단계를 거치면 실제 현장 배포 전에 품질을 한층 끌어올릴 수 있습니다.

以上的 로드맵과 전략을 통해, 우리 스터디 팀은 **임베딩 환경용 LLM 개발**이라는 도전적 목표를 단계적으로 달성할 수 있을 것입니다. 중요한 것은 **협업을 통해 지속적으로 학습**하고 **문제 해결을 즐기는 태도**입니다. 각 단계에서 얻은 지식을 토대로 다음 단계를 수월하게 밟아나가며, 프로젝트 완료 시점에는 팀원 모두가 LLM 모델 경량화와 엣지 배포에 대한 전문성을 한층 높이게 되길 기대합니다.

**참고한 자료 및 출처:** 본 문서에서는 최신 경량화 기술 동향 <sup>5</sup> <sup>2</sup>, Kaggle 경진대회 정보 <sup>19</sup> <sup>21</sup> 등을 인용하여 내용을 구성하였습니다. 자세한 내용은 각 출처를 통해 살펴볼 수 있습니다. 현 시점(2025년)에도 관련 기술과 사례는 빠르게 발전 중이므로, 꾸준한 정보 업데이트와 추가 실험을 병행하면 더욱 탄탄한 프로젝트 수행이 가능할 것입니다. ●

<sup>1</sup> <sup>4</sup> <sup>6</sup> <sup>8</sup> <sup>11</sup> <sup>12</sup> <sup>13</sup> <sup>14</sup> <sup>15</sup> <sup>16</sup> <sup>17</sup> <sup>18</sup> Deploying LLMs in Low-Resource Environments: Edge AI, KV Caching, and Model Quantization | by Tommy Adeliyi | Medium

<https://medium.com/@tommyadeliyi/deploying-llms-in-low-resource-environments-edge-ai-kv-caching-and-model-quantization-a119c2df7716>

<sup>2</sup> <sup>5</sup> <sup>7</sup> Edge deployment of LLMs and ML models: A review

<https://www.rohan-paul.com/p/edge-deployment-of-llms-and-ml-models>

<sup>3</sup> <sup>9</sup> <sup>10</sup> <sup>19</sup> <sup>20</sup> Kaggle Competition Report: LLM Science Exam | Hippocampus's Garden

[https://hippocampus-garden.com/kaggle\\_llm/](https://hippocampus-garden.com/kaggle_llm/)

<sup>21</sup> <sup>22</sup> <sup>23</sup> Trying to solve a Kaggle Competition: LLM — Detect AI Generated Text (Part 1) | by Odilbek Tokhirov | Medium

<https://medium.com/@odil.tokhirov/trying-to-solve-a-kaggle-competition-llm-detect-ai-generated-text-part-1-c53b1b087795>

<sup>24</sup> Text Authorship Classification: An Overview of the Kaggle Competition 'H2O Predict the LLM' | by Samvelkoch | Medium

<https://medium.com/@samvelkoch/text-authorship-classification-an-overview-of-the-kaggle-competition-h2o-predict-the-llm-bbbf874972cf>

25 Kaggle Solutions

<https://farid.one/kaggle-solutions/>