

# Open and Efficient Foundation Language Models

## 1. 논문의 핵심 주제

- **LLaMA**는 Meta에서 공개한 언어모델 시리즈로, 7B~65B 파라미터 규모로 구성돼 있어.
- 공개 데이터로만 학습해도 최상위급 성능이 가능함을 보여주는 게 목적이야.

---

## 2. LLaMA 모델의 장점 🚀

- **효율성:**
  - LLaMA-13B는 GPT-3 (175B)를 여러 벤치마크에서 능가하면서도 10배 이상 작음.
  - LLaMA-65B는 PaLM-540B와 Chinchilla-70B 같은 초대형 모델과도 경쟁할 수준.
- **개방성:**
  - 다른 초대형 모델과 달리, 오직 공개된 데이터만 사용 → 연구 커뮤니티에 모델을 오픈소스로 제공할 수 있음.
- **경제성:**
  - 작은 모델을 더 많은 데이터로 오래 학습하는 방식이 효율적인 성능을 가져옴.
  - 추론(Inference) 비용을 절약할 수 있어 서비스에 활용하기 좋음.

---

## 3. 데이터 구성 (총 1.4조 토큰 📦)

데이터셋	비율 특징
CommonCrawl	67% 인터넷 크롤링 데이터, 고품질 페이지만 추출
C4	15% 전처리된 웹 크롤링 데이터 (Google 제작)
GitHub	4.5% 오픈소스 코드
Wikipedia	4.5% 다국어 위키백과
Gutenberg & Books3	4.5% 무료 배포 가능한 책 데이터
ArXiv	2.5% 과학 논문 데이터

데이터셋	비율	특징
Stack Exchange	2%	Q&A 형식 데이터

#### 4. 모델 구조 🧠 (Transformer 기반)

- **Transformer 개선점:**
  - Pre-normalization (안정성 향상)
  - SwiGLU 활성화 함수 (성능 향상)
  - Rotary positional embeddings (더 나은 위치 정보 활용)

##### 모델 크기 차원 헤드 수 레이어 수 학습률 학습 토큰 수

7B	4096	32	32	3.0e-4	1.0T
13B	5120	40	40	3.0e-4	1.0T
33B	6656	52	60	1.5e-4	1.4T
65B	8192	64	80	1.5e-4	1.4T

- **Optimizer:** AdamW, Cosine learning rate 사용
- **학습 속도 최적화:** GPU 간 통신 병렬화, 메모리 최적화 등으로 2048개의 A100 GPU에서 65B 모델을 21일 만에 학습 완료함.

#### 5. 주요 결과 및 성능 비교 📊

- **상식 추론 (Commonsense Reasoning):**
  - LLaMA-65B는 Chinchilla-70B, PaLM-540B를 능가하거나 비슷한 성능
  - LLaMA-13B는 GPT-3를 대부분 능가
- **질의응답 (Closed-book QA):**
  - LLaMA-65B가 zero-shot, few-shot 설정에서 모두 최상급 성능 기록 (TriviaQA 등)
- **읽기 이해 (Reading Comprehension):**
  - LLaMA-13B는 GPT-3를 능가, 65B는 PaLM-540B와 경쟁 가능

- **수리적 추론 (Mathematical Reasoning):**
    - LLaMA-65B가 fine-tune 없는 상태에서도 Minerva-62B와 같은 수학 전문 모델을 넘는 성능 보임 (GSM8k 기준)
  - **코드 생성 (Code Generation):**
    - 비슷한 크기의 다른 언어모델보다 더 우수한 성능 (HumanEval, MBPP 벤치마크 기준)
  - **다중 작업 언어 이해 (MMLU):**
    - LLaMA-65B는 Chinchilla, PaLM과 비교하여 약간 낮은 성능 (책 데이터 양의 부족 때문)
- 

## 6. 모델의 한계와 이슈 !

- **독성 (Toxicity):**
    - 모델 크기가 커질수록 독성 발언 생성 가능성 증가
  - **편향성 (Bias):**
    - 특정 성별, 종교, 직업 등에 대한 편향 존재함 (CrowS-Pairs, WinoGender 기준)
  - **허구 정보 생성 (Truthfulness):**
    - 잘못된 정보 생성 가능성 존재 (TruthfulQA에서 일부 개선된 성능 보이지만 한계 존재)
- 

## 7. 탄소 배출량 🌱

- LLaMA-65B의 학습 과정에서 약 **173톤의 탄소**를 배출.
  - 오픈소스로 공개되었으므로, 재학습으로 인한 추가 탄소 배출을 줄이는 효과 기대.
- 

## 8. 결론 및 향후 방향 🚀

- 개방형 데이터로도 최상급 모델 학습이 가능하다는 점을 증명.
- LLaMA 시리즈를 공개하여 연구자들의 언어모델 연구를 가속화 기대.
- 향후 더 큰 데이터와 모델을 공개하고, Instruction finetuning(지시어 기반 미세조정)을 발

전시킬 예정.

---

#### 발표 핵심 키워드

- LLaMA
- 개방형 데이터
- 모델 효율성 & 경제성
- 초대형 언어모델과의 경쟁력

#### 결론 (논문의 핵심 take-away)

- 초대형 언어모델 성능 향상의 핵심은 모델 크기 확대뿐 아니라, 더 많은 데이터로 작은 모델을 장기간 학습시키는 전략.
- 공개 데이터로만 학습된 고성능 오픈소스 언어모델(LLaMA)을 통해 연구 커뮤니티의 발전 촉진 가능.