



# LLaMA & TCN

날짜	@2025년 3월 20일
분야	분석
주차	과제

## LLaMA

### 1. 논문 개요

논문 제목: LLaMA: Open and Efficient Foundation Language Models

저자: Hugo Touvron et al. (Meta AI)

발표 연도: 2023년

출처: arXiv

Meta AI에서 발표한 LLaMA(Large Language Model Meta AI) 시리즈는 기존 대형 언어 모델 대비 효율적인 학습과 추론 성능을 목표로 한다. LLaMA는 작은 모델에서 더 많은 데이터를 학습하는 것이 성능 향상에 효과적이라는 최근 연구 결과를 반영하여 개발되었으며, 오픈소스로 공개되어 연구 커뮤니티에서 활용할 수 있도록 설계되었다.

### 2. 연구의 주요 기여

- 효율적인 학습 전략:** 작은 모델이 더 많은 데이터로 학습하는 것이 compute budget 대비 최적의 성능을 낼 수 있음을 보였다.
- 다양한 모델 크기 제공:** 7B, 13B, 30B, 65B 파라미터 규모의 모델을 제공하여 다양한 사용 사례에 적용 가능.
- 공개 데이터셋 활용:** Chinchilla, PaLM, GPT-3와 달리 오직 공개적으로 접근 가능한 데이터만을 사용하여 학습.
- 추론 비용 최적화:** 작은 모델을 장기간 학습하여 inference 비용을 절감하는 방식을 채택.
- 오픈 액세스 지원:** 연구 커뮤니티에서 접근 가능하도록 공개하여 확장성과 응용성을 높임.

### 3. 모델 구조 및 학습 방식

#### 3-1. 학습 데이터

- 데이터 출처:** Wikipedia, Common Crawl, C4, GitHub, Books, arXiv, Stack Exchange 등 다양한 공개 데이터셋 활용.
- 토큰화 방법:** Byte Pair Encoding(BPE) 사용, 숫자는 개별 문자로 변환하여 처리.
- 총 학습 데이터 규모:** 1.4T 토큰 사용, 일부 데이터(Wikipedia, Books)는 2 epoch 수행.

#### 3-2. 모델 아키텍처

- Transformer 기반 모델**
- Pre-normalization:** GPT-3에서 사용된 방법을 채택하여 학습 안정성 개선.
- SwiGLU 활성화 함수:** ReLU 대신 SwiGLU를 적용하여 성능 향상.
- Rotary Positional Embedding(RoPE):** GPT-Neo와 동일한 방식으로 positional encoding 최적화.
- AdamW Optimizer 사용:** Cosine learning rate 스케줄 적용, weight decay 0.1, gradient clipping 1.0.

#### 3-3. 효율적인 구현

- Casual multi-head attention 사용:** 메모리 사용량 감소 및 연산량 절감.
- Backward pass 최적화:** activation 재계산을 최소화하여 학습 속도 향상.
- 모델 및 시퀀스 병렬화:** 메모리 효율성을 극대화.

## 4. 실험 결과

### 4-1. Zero-shot 및 Few-shot 평가

- 벤치마크: GPT-3, Gopher, Chinchilla, PaLM, OPT, GPT-J, GPT-Neo와 비교 수행.
- 평가 방식
  - Zero-shot: Task 설명 및 test example만 제공.
  - Few-shot: Few-shot learning 예제 제공 후 평가.

### 4-2. MMLU 벤치마크 결과

- 5-shot setting에서 평가 수행.
- LLaMA-65B는 Chinchilla-70B, PaLM-540B와 비교하여 약간 낮은 성능을 보였으나, 학습 데이터에서 books 및 academic paper의 비중이 적었던 것이 원인으로 분석됨.
- Perplexity와 성능 간 상관 관계 확인, training loss가 감소함에 따라 성능 향상 관찰됨.

## 5. Instruction Fine-tuning

- 간단한 instruction fine-tuning만으로도 MMLU 성능 개선 가능.
- Instruction-tuned 모델 LLaMA-I는 기존 instruction-tuned 모델(OPT-IML, Flan-PaLM)과 비교하여 높은 성능 달성(68.9%).
- Fine-tuning을 통한 추가적인 task 적응 가능성 확인.

## 6. 결론

LLaMA는 공개적으로 접근 가능한 데이터만을 사용하여 기존 대형 언어 모델(GPT-3, PaLM, Chinchilla)과 경쟁할 수 있는 성능을 보여주었다. 특히, LLaMA-13B는 GPT-3보다 10배 작은 모델임에도 불구하고 동등한 성능을 달성했으며, LLaMA-65B는 Chinchilla-70B 및 PaLM-540B와 유사한 성능을 보였다. 또한, instruction fine-tuning을 통해 추가적인 성능 향상이 가능함을 입증하였다. 무엇보다 LLaMA는 연구 커뮤니티에 공개됨으로써 LLM 연구 및 응용을 보다 용이하게 만들었다.

## TCN

### 1. 논문 개요

본 논문은 시퀀스 모델링에서 Temporal Convolutional Network(TCN)과 Recurrent Neural Networks(RNN) 계열 모델(LSTM, GRU)을 비교 분석한 연구입니다. 저자들은 다양한 시퀀스 예측 및 시계열 분석 태스크에서 TCN과 RNN의 성능을 실험적으로 평가하였습니다.

### 2. 연구 목적

시퀀스 모델링에서 RNN 계열이 전통적으로 많이 사용되었지만, CNN 기반의 접근법인 TCN이 RNN보다 장기 의존성(Long-term dependency)을 더 잘 학습하고, 병렬 처리가 가능하며, 학습 안정성이 뛰어나다는 점을 강조합니다. 본 연구에서는 다양한 벤치마크 테스트를 통해 TCN과 RNN 계열의 모델을 정량적으로 비교하였습니다.

### 3. TCN(Temporal Convolutional Network)의 특징

TCN은 기본적으로 1D CNN을 기반으로 시퀀스를 모델링하는 방식으로 다음과 같은 특징을 가집니다.

#### 3.1 Causal Convolution

- 시계열 데이터의 시간 순서를 유지하기 위해 Causal Convolution을 사용합니다.
- 즉, 출력 시점  $t$ 는  $t$  이전의 입력 값에만 영향을 받으며, 미래 데이터를 참조하지 않습니다.

#### 3.2 Dilated Convolution

- Dilated Convolution을 사용하여 receptive field를 넓혀, 장기 의존성을 효과적으로 학습할 수 있습니다.
- 일반적인 CNN보다 적은 계층 수로 더 넓은 시간 범위를 커버할 수 있습니다.

### 3.3 Residual Connections

- Residual Connection을 도입하여 깊은 네트워크에서도 **gradient vanishing 문제**를 방지하고, 학습 안정성을 높입니다.

### 3.4 Fully Convolutional Architecture

- RNN과 달리 **순차적 연산이 아니라 병렬 연산이 가능**하여, 학습 속도가 빠릅니다.

## 4. 실험 및 성능 비교

논문에서는 다양한 시퀀스 태스크에서 TCN과 RNN(LSTM, GRU)을 비교 평가하였습니다.

### 4.1 벤치마크 태스크

- **Adding Problem**: 장기 의존성을 평가하는 문제
- **Sequential MNIST/PERMUTED MNIST**: 이미지 데이터에서 시퀀스 예측 문제
- **Copy Memory Task**: 장기 기억 유지 능력 평가
- **Char-level Language Modeling**: 언어 모델링 태스크

### 4.2 실험 결과

- **Adding Problem, Copy Memory Task** 등 장기 의존성이 중요한 문제에서는 **TCN이 RNN보다 우수한 성능을 보임**
- **Sequential MNIST/PERMUTED MNIST** 실험에서도 TCN이 더 높은 정확도를 기록
- **Char-level Language Modeling**에서는 RNN(LSTM, GRU)이 일부 더 좋은 성능을 보였으나, TCN도 경쟁력 있는 결과를 냄
- TCN은 **병렬 처리 가능**, 학습 속도가 빠르고 **RNN보다 안정적인 성능**을 보임

## 5. 결론 및 시사점

- TCN은 시퀀스 모델링에서 RNN보다 **장기 의존성을 더 효과적으로 학습**할 수 있으며, **병렬 연산이 가능**하여 연산 속도에서 이점이 있음.
- Residual Connections와 Dilated Convolution을 활용하여 **학습 안정성**이 뛰어남.
- 그러나 언어 모델링과 같은 특정 문제에서는 여전히 RNN 계열이 유리한 경우도 있음.
- 결론적으로, RNN을 대체할 수 있는 **강력한 대안으로 TCN을 고려할 가치가 있음**.

## 6. 개인적인 평가 및 향후 연구 방향

- TCN은 **시계열 예측, 음성 처리, 금융 데이터 분석** 등 장기 의존성이 중요한 문제에서 강력한 대안이 될 수 있음.
- 그러나 NLP 분야에서 Transformer 계열 모델이 등장하면서, TCN이 얼마나 경쟁력을 가질지는 추가 연구가 필요함.
- 또한, TCN의 **메모리 사용량 및 계산량 최적화**에 대한 연구도 진행될 필요가 있음.

향후 연구 방향:

1. **TCN과 Transformer를 결합한 모델 연구**
2. **TCN의 구조 최적화를 통한 계산 비용 절감**
3. **다양한 실제 애플리케이션에서 TCN 적용 사례 분석**