



# LLaMA: Open and Efficient Foundation Language Models

김 수

# 목차

- Introduction
- Approach
- Main results
- Bias, Toxicity and Misinformation
- Conclusion

# Introduction

- Large Language Models(LLM)
- 모델 크기가 작아도(7B~64B) 더 많은 토큰을 훈련하여 최상의 성능 달성
- 공개적으로 이용 가능한 데이터만 사용

# Approach

- 이전 연구에서 묘사된 방법과 유사함
- Chinchila scaling law에서 영감을 받음
- 표준 Optimizer 사용하여 대량의 텍스트 데이터에 대한 큰 Transformer 훈련함

# Approach

- Data

- 공개적으로 이용 가능한 Dataset을 사용하여 전처리함

Dataset	Sampling prop.	Epochs	Disk size
CommonCrawl	67.0%	1.10	3.3 TB
C4	15.0%	1.06	783 GB
Github	4.5%	0.64	328 GB
Wikipedia	4.5%	2.45	83 GB
Books	4.5%	2.23	85 GB
ArXiv	2.5%	1.06	92 GB
StackExchange	2.0%	1.03	78 GB

Table 1: **Pre-training data.** Data mixtures used for pre-training, for each subset we list the sampling proportion, number of epochs performed on the subset when training on 1.4T tokens, and disk size. The pre-training runs on 1T tokens have the same sampling proportion.

# Approach

- Architecture
  - Pre-normalization [GPT3]
  - SwiGLU activation function [PaLM]
  - Rotary positional embeddings[GPTNeo]

# Approach

- Efficient implementation
  - The casual multi-head attention
  - Reduced the amount of activations
  - Reduce the memory usage of the model

# Main results

- Zero-shot
  - 작업의 텍스트 설명과 테스트 예제 제공
  - open-ended generation 사용하여 답 제공하거나 제안된 답안을 순위로 나열함.
- Few-shot
  - 몇 가지 예제(1~64개)와 테스트 예제 제공
  - 이 텍스트를 입력으로 받아 답을 생성하거나 다양한 옵션을 순위로 나열함



# Main results

- Common Sence Reasoning
  - 다지 선다형 등

		BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA
GPT-3	175B	60.5	81.0	-	78.9	70.2	68.8	51.4	57.6
Gopher	280B	79.3	81.8	50.6	79.2	70.1	-	-	-
Chinchilla	70B	83.7	81.8	51.3	80.8	74.9	-	-	-
PaLM	62B	84.8	80.5	-	79.7	77.0	75.2	52.5	50.4
PaLM-cont	62B	83.9	81.4	-	80.6	77.0	-	-	-
PaLM	540B	<b>88.0</b>	82.3	-	83.4	<b>81.1</b>	76.6	53.0	53.4
LLaMA	7B	76.5	79.8	48.9	76.1	70.1	72.8	47.6	57.2
	13B	78.1	80.1	50.4	79.2	73.0	74.8	52.7	56.4
	33B	83.1	82.3	50.4	82.8	76.0	<b>80.0</b>	<b>57.8</b>	58.6
	65B	85.3	<b>82.8</b>	<b>52.3</b>	<b>84.2</b>	77.0	78.9	56.0	<b>60.2</b>

Table 3: Zero-shot performance on Common Sense Reasoning tasks.

# Main results

- Closed-book Question Answering
  - 외부 지식을 활용하지 않고 사전에 학습한 지식을 이용하여 정답을 도출함

		0-shot	1-shot	5-shot	64-shot
GPT-3	175B	14.6	23.0	-	29.9
Gopher	280B	10.1	-	24.5	28.2
Chinchilla	70B	16.6	-	31.5	35.5
PaLM	8B	8.4	10.6	-	14.6
	62B	18.1	26.5	-	27.6
	540B	21.2	29.3	-	39.6
LLaMA	7B	16.8	18.7	22.0	26.1
	13B	20.1	23.4	28.1	31.9
	33B	<b>24.9</b>	28.3	32.9	36.0
	65B	23.8	<b>31.0</b>	<b>35.0</b>	<b>39.9</b>

Table 4: **NaturalQuestions**. Exact match performance.

		0-shot	1-shot	5-shot	64-shot
Gopher	280B	43.5	-	57.0	57.2
Chinchilla	70B	55.4	-	64.1	64.6
LLaMA	7B	50.0	53.4	56.3	57.6
	13B	56.6	60.5	63.1	64.0
	33B	65.1	67.9	69.9	70.4
	65B	<b>68.2</b>	<b>71.6</b>	<b>72.6</b>	<b>73.0</b>

Table 5: **TriviaQA**. Zero-shot and few-shot exact match performance on the filtered dev set.

# Main results

- Reading Comprehension
  - 객관식 문제 풀기

		RACE-middle	RACE-high
GPT-3	175B	58.4	45.5
PaLM	8B	57.9	42.3
	62B	64.3	47.5
	540B	<b>68.1</b>	49.1
LLaMA	7B	61.1	46.9
	13B	61.6	47.2
	33B	64.1	48.3
	65B	67.9	<b>51.6</b>

Table 6: **Reading Comprehension.** Zero-shot accuracy.

# Main results

- Mathematical Reasoning
  - 수학 문제(문제, 풀이과정, 정답), 정답 맞추면 됨

		MATH +maj1@k		GSM8k +maj1@k	
PaLM	8B	1.5	-	4.1	-
	62B	4.4	-	33.0	-
	540B	8.8	-	56.5	-
Minerva	8B	14.1	25.4	16.2	28.4
	62B	27.6	43.4	52.4	68.5
	540B	<b>33.6</b>	<b>50.3</b>	<b>68.5</b>	<b>78.5</b>
LLaMA	7B	2.9	6.9	11.0	18.1
	13B	3.9	8.8	17.8	29.3
	33B	7.1	15.2	35.6	53.1
	65B	10.6	20.5	50.9	69.7

Table 7: **Model performance on quantitative reasoning datasets.** For majority voting, we use the same setup as Minerva, with  $k = 256$  samples for MATH and  $k = 100$  for GSM8k (Minerva 540B uses  $k = 64$  for MATH and  $k = 40$  for GSM8k). LLaMA-65B outperforms Minerva 62B on GSM8k, although it has not been fine-tuned on mathematical data.

# Main results

- Code Generation
  - 주어진 자연어 문장을 바탕으로 코드 생성함
  - Unit test case 들을 통과해야 함

pass@	Params	HumanEval		MBPP	
		@1	@100	@1	@80
LaMDA	137B	14.0	47.3	14.8	62.4
PaLM	8B	3.6*	18.7*	5.0*	35.7*
PaLM	62B	15.9	46.3*	21.4	63.2*
PaLM-cont	62B	23.7	-	31.2	-
PaLM	540B	<b>26.2</b>	76.2	36.8	75.0
LLaMA	7B	10.5	36.5	17.7	56.2
	13B	15.8	52.5	22.0	64.0
	33B	21.7	70.7	30.2	73.4
	65B	23.7	<b>79.3</b>	<b>37.7</b>	<b>76.8</b>

Table 8: **Model performance for code generation.** We report the pass@ score on HumanEval and MBPP. HumanEval generations are done in zero-shot and MBPP with 3-shot prompts similar to [Austin et al. \(2021\)](#). The values marked with \* are read from figures in [Chowdhery et al. \(2022\)](#).

# Main results

- MMLU: Massive Multitask Language Understanding

- 다양한 지식 분야의 객관식 문제

		Humanities	STEM	Social Sciences	Other	Average
GPT-NeoX	20B	29.8	34.9	33.7	37.7	33.6
GPT-3	175B	40.8	36.7	50.4	48.8	43.9
Gopher	280B	56.2	47.4	71.9	66.1	60.0
Chinchilla	70B	63.6	54.9	79.3	<b>73.9</b>	67.5
PaLM	8B	25.6	23.8	24.1	27.8	25.4
	62B	59.5	41.9	62.7	55.8	53.7
	540B	<b>77.0</b>	<b>55.6</b>	<b>81.0</b>	69.6	<b>69.3</b>
LLaMA	7B	34.0	30.5	38.3	38.1	35.1
	13B	45.0	35.8	53.8	53.3	46.9
	33B	55.8	46.0	66.7	63.4	57.8
	65B	61.8	51.7	72.9	67.4	63.4

Table 9: **Massive Multitask Language Understanding (MMLU)**. Five-shot accuracy.

OPT	30B	26.1
GLM	120B	44.8
PaLM	62B	55.1
PaLM-cont	62B	62.8
Chinchilla	70B	67.5
LLaMA	65B	63.4
OPT-IML-Max	30B	43.2
Flan-T5-XXL	11B	55.1
Flan-PaLM	62B	59.6
Flan-PaLM-cont	62B	66.1
LLaMA-I	65B	<b>68.9</b>

Table 10: **Instruction finetuning – MMLU (5-shot)**. Comparison of models of moderate size with and without instruction finetuning on MMLU.

# Bias, Toxicity and Misinformation

- RealToxicityPrompts
  - 목적: 독성을 가진 언어가 등장하는가?

		Basic	Respectful
LLaMA	7B	0.106	0.081
	13B	0.104	0.095
	33B	0.107	0.087
	65B	0.128	0.141

Table 11: **RealToxicityPrompts**. We run a greedy decoder on the 100k prompts from this benchmark. The “respectful” versions are prompts starting with “Complete the following sentence in a polite, respectful, and unbiased manner:”, and “Basic” is without it. Scores were obtained using the PerplexityAPI, with higher score indicating more toxic generations.

# Bias, Toxicity and Misinformation

- CrowS-Pairs

- 편향성 측정
- 높을 수록 많은 편향을 가짐

	LLaMA	GPT3	OPT
Gender	70.6	<b>62.6</b>	65.7
Religion	79.0	73.3	<b>68.6</b>
Race/Color	<b>57.0</b>	64.7	68.6
Sexual orientation	81.0	<b>76.2</b>	78.6
Age	70.1	<b>64.4</b>	67.8
Nationality	64.2	<b>61.6</b>	62.9
Disability	<b>66.7</b>	76.7	76.7
Physical appearance	77.8	<b>74.6</b>	76.2
Socioeconomic status	<b>71.5</b>	73.8	76.2
Average	<b>66.6</b>	67.2	69.5

Table 12: **CrowS-Pairs.** We compare the level of biases contained in LLaMA-65B with OPT-175B and GPT3-175B. Higher score indicates higher bias.



# Bias, Toxicity and Misinformation

- WinoGender

- 성별 편향

	7B	13B	33B	65B
All	66.0	64.7	69.0	77.5
her/her/she	65.0	66.7	66.7	78.8
his/him/he	60.8	62.5	62.1	72.1
their/them/someone	72.1	65.0	78.3	81.7
her/her/she ( <i>gotcha</i> )	64.2	65.8	61.7	75.0
his/him/he ( <i>gotcha</i> )	55.0	55.8	55.8	63.3

Table 13: **WinoGender**. Co-reference resolution accuracy for the LLaMA models, for different pronouns (“her/her/she” and “his/him/he”). We observe that our models obtain better performance on “their/them/someone” pronouns than on “her/her/she” and “his/him/he”, which is likely indicative of biases.

# Bias, Toxicity and Misinformation

## • TruthfulQA

- 어떤 주장에 대해 참을 얘기하는가?
- 거짓 정보 또는 허위 주장을 생성하는가?

		Truthful	Truthful*Inf
GPT-3	1.3B	0.31	0.19
	6B	0.22	0.19
	175B	0.28	0.25
LLaMA	7B	0.33	0.29
	13B	0.47	0.41
	33B	0.52	0.48
	65B	0.57	0.53

Table 14: **TruthfulQA**. We report the fraction of truthful and truthful\*informative answers, as scored by specially trained models via the OpenAI API. We follow the QA prompt style used in [Ouyang et al. \(2022\)](#), and report the performance of GPT-3 from the same paper.

# Conclusion

- GPT-3 보다 성능이 좋으면서 모델 크기가 10배 이상 작음
- Chinchilla-70B and PaLM-540B와 경쟁 가능한 정도의 성능을 보임