



Mistral 7B

김 수

목차

- Abstract
- Introduction
- Architectural details
- Results
- Instruction Finetuning
- Adding guardrails for front-facing applications
- Conclusion

Abstract

- Mistral 7B
 - 7Billion parameter language model
 - Outperforms in reasoning, mathematics, and code generation
 - Llama2(13B model)
 - Llama(34B model)

Introduction

- Mistral 7B
 - 효율적인 결과 유지하면서 높은 성능 제공하는 모델
 - 성능
 - Llama2-13B : 모든 테스트된 벤치마크 부분 성능 증가
 - Llama-34B : 수학, 코드 생성 부분 성능 증가

Introduction

- Mistral 7B
 - Grouped-Query Attention(GQA)
 - Inference time speed 가속화
 - 디코딩 중 메모리의 필요 요구량을 줄임
 - 더 높은 배치 크기 허용하여 **처리량을 향상**시킴
 - Sliding Window Attention(SWA)
 - 계산 비용 줄이면서 긴 시퀀스 길이를 **효율적**으로 처리
 - LLA에서 한계를 완화함

Architectural details

- Main parameters

| Parameter | Value |
|-------------|-------|
| dim | 4096 |
| n_layers | 32 |
| head_dim | 128 |
| hidden_dim | 14336 |
| n_heads | 32 |
| n_kv_heads | 8 |
| window_size | 4096 |
| context_len | 8192 |
| vocab_size | 32000 |

Table 1: Model architecture.

Architectural details

- Sliding Window Attention

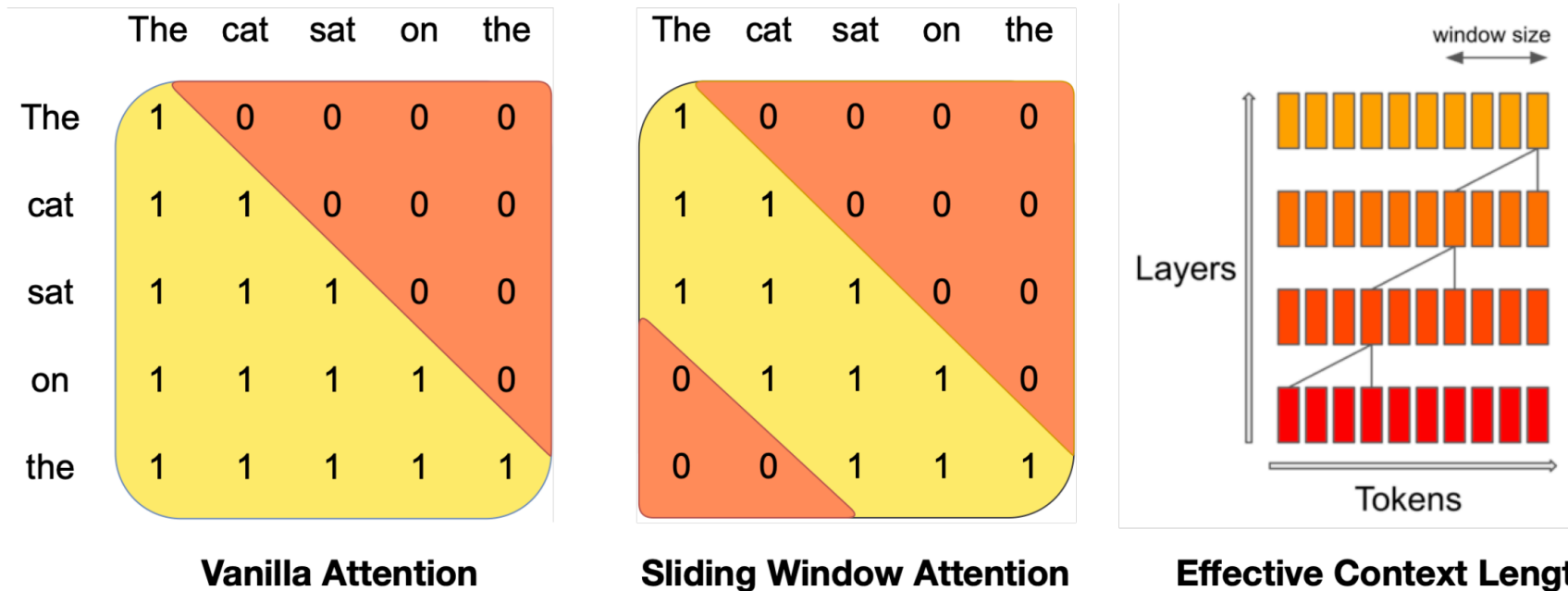


Figure 1: Sliding Window Attention. The number of operations in vanilla attention is quadratic in the sequence length, and the memory increases linearly with the number of tokens. At inference time, this incurs higher latency and smaller throughput due to reduced cache availability. To alleviate this issue, we use sliding window attention: each token can attend to at most W tokens from the previous layer (here, $W = 3$). Note that tokens outside the sliding window still influence next word prediction. At each attention layer, information can move forward by W tokens. Hence, after k attention layers, information can move forward by up to $k \times W$ tokens.

Architectural details

- Rolling Buffer Cache

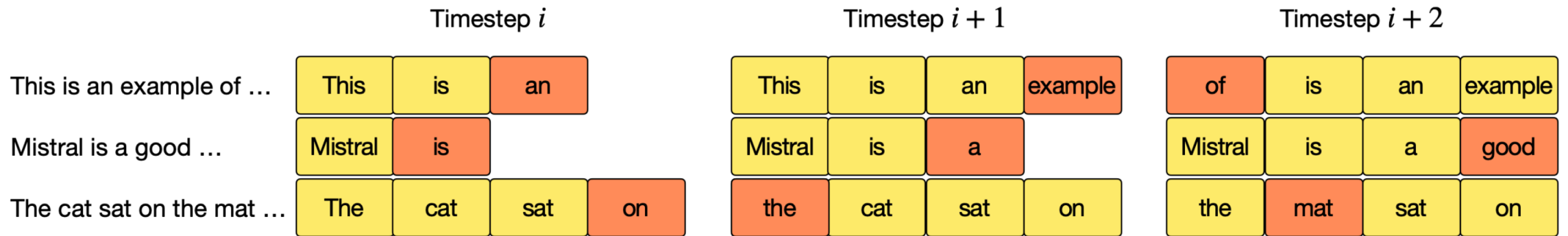


Figure 2: Rolling buffer cache. The cache has a fixed size of $W = 4$. Keys and values for position i are stored in position $i \bmod W$ of the cache. When the position i is larger than W , past values in the cache are overwritten. The hidden state corresponding to the latest generated tokens are colored in orange.

Architectural details

- Pre-fill and Chunking

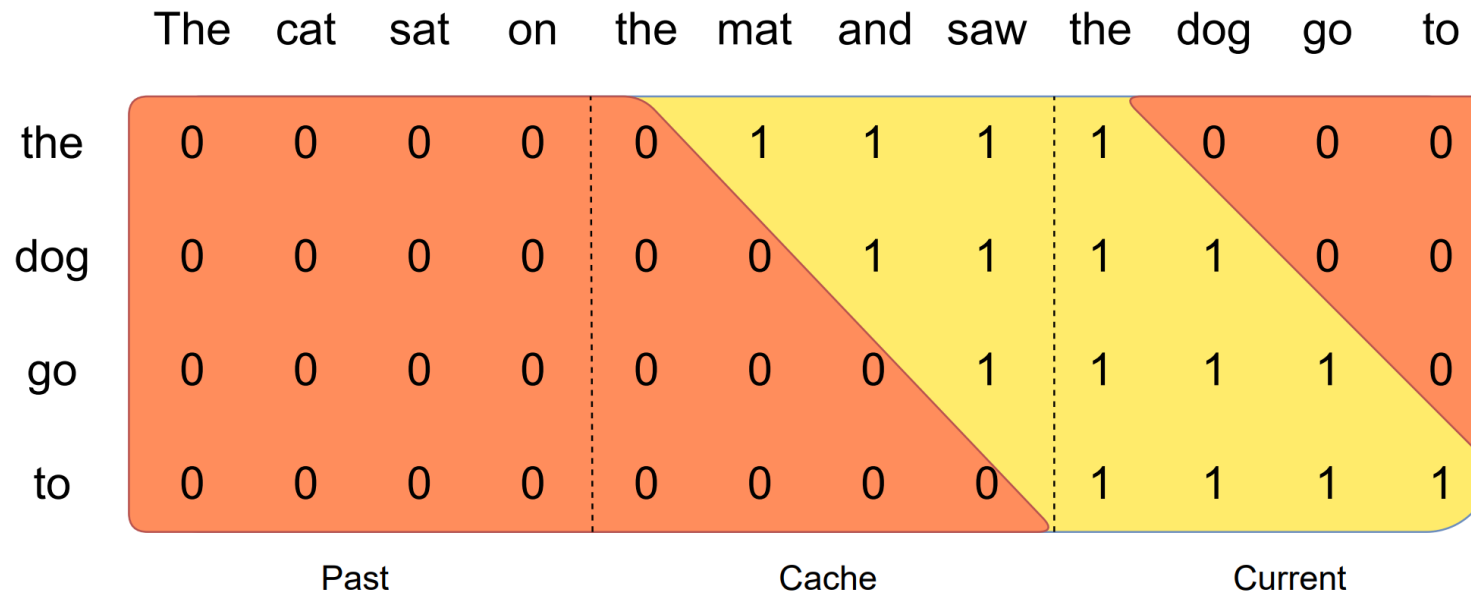


Figure 3: Pre-fill and chunking. During pre-fill of the cache, long sequences are chunked to limit memory usage. We process a sequence in three chunks, "The cat sat on", "the mat and saw", "the dog go to". The figure shows what happens for the third chunk ("the dog go to"): it attends itself using a causal mask (rightmost block), attends the cache using a sliding window (center block), and does not attend to past tokens as they are outside of the sliding window (left block).

Results

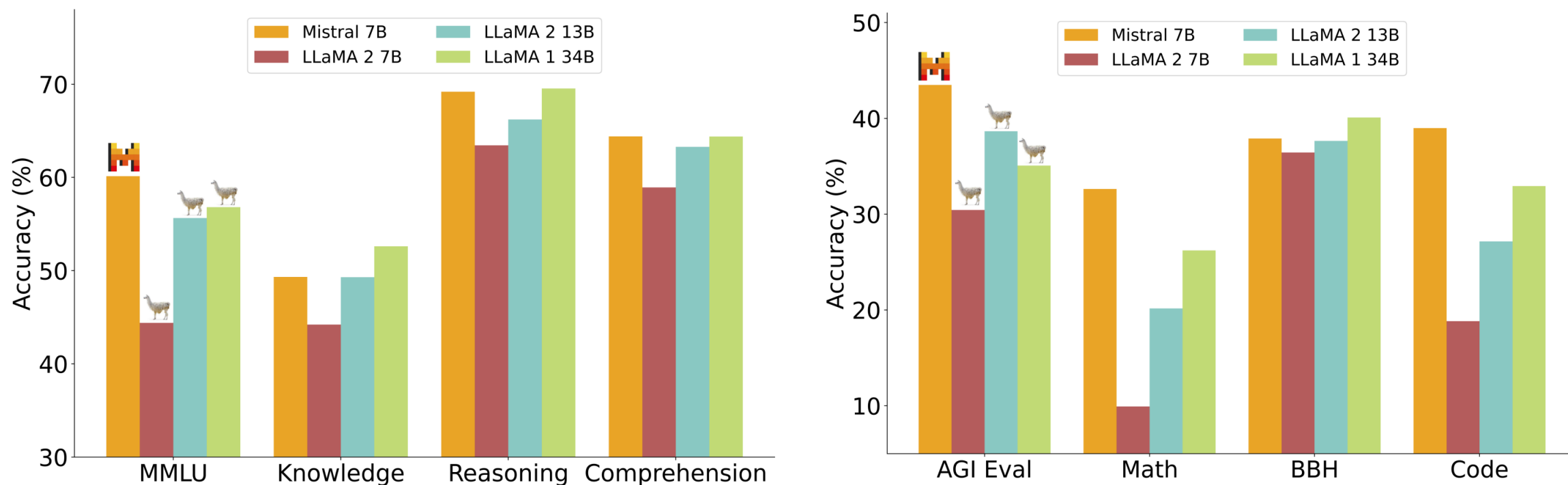


Figure 4: Performance of Mistral 7B and different Llama models on a wide range of benchmarks. All models were re-evaluated on all metrics with our evaluation pipeline for accurate comparison. Mistral 7B significantly outperforms Llama 2 7B and Llama 2 13B on all benchmarks. It is also vastly superior to Llama 1 34B in mathematics, code generation, and reasoning benchmarks.

Results

Popular aggregated results

Commonsense Reasoning

World Knowledge

Code

Math

| Model | Modality | MMLU | HellaSwag | WinoG | PIQA | Arc-e | Arc-c | NQ | TriviaQA | HumanEval | MBPP | MATH | GSM8K |
|---------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| LLaMA 2 7B | Pretrained | 44.4% | 77.1% | 69.5% | 77.9% | 68.7% | 43.2% | 24.7% | 63.8% | 11.6% | 26.1% | 3.9% | 16.0% |
| LLaMA 2 13B | Pretrained | 55.6% | 80.7% | 72.9% | 80.8% | 75.2% | 48.8% | 29.0% | 69.6% | 18.9% | 35.4% | 6.0% | 34.3% |
| Code-Llama 7B | Finetuned | 36.9% | 62.9% | 62.3% | 72.8% | 59.4% | 34.5% | 11.0% | 34.9% | 31.1% | 52.5% | 5.2% | 20.8% |
| Mistral 7B | Pretrained | 60.1% | 81.3% | 75.3% | 83.0% | 80.0% | 55.5% | 28.8% | 69.9% | 30.5% | 47.5% | 13.1% | 52.2% |

Table 2: Comparison of Mistral 7B with Llama. Mistral 7B outperforms Llama 2 13B on all metrics, and approaches the code performance of Code-Llama 7B without sacrificing performance on non-code benchmarks.

Results

- Knowledge 1.9x의 낮은 압축률

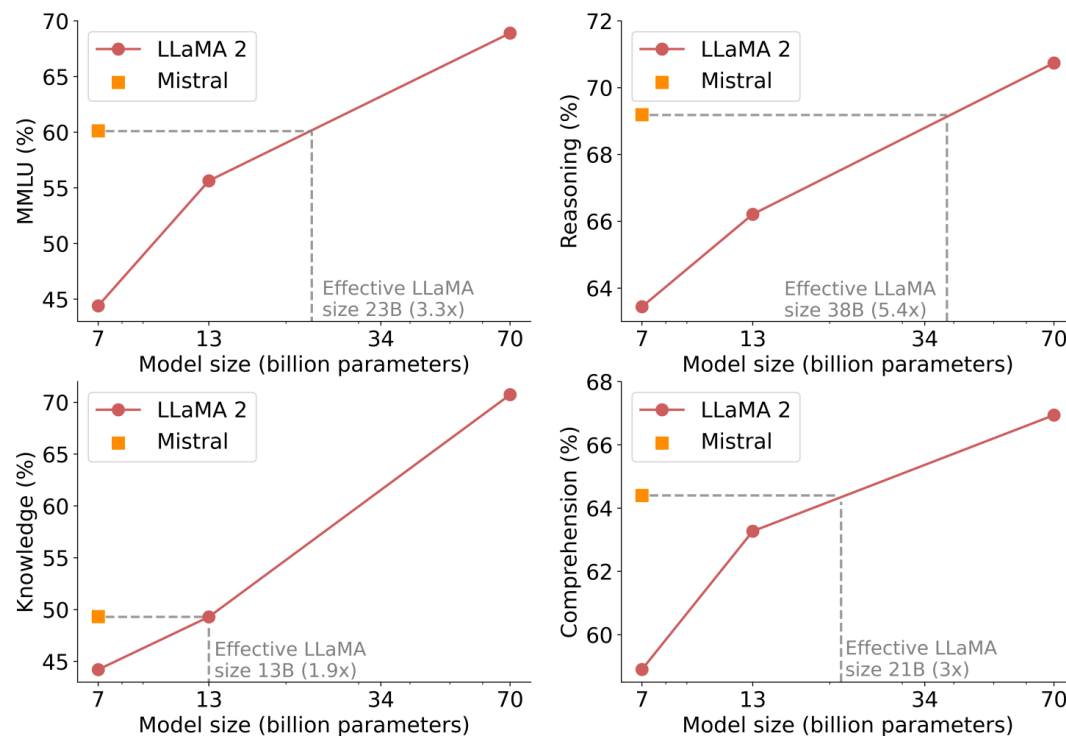


Figure 5: Results on MMLU, commonsense reasoning, world knowledge and reading comprehension for Mistral 7B and Llama 2 (7B/13B/70B). Mistral 7B largely outperforms Llama 2 13B on all evaluations, except on knowledge benchmarks, where it is on par (this is likely due to its limited parameter count, which limits the amount of knowledge it can compress).

Instruction Finetuning

- Fine-tuned it on instruction datasets

| Model | Chatbot Arena ELO Rating | MT Bench |
|----------------------------|-----------------------------|----------------------|
| WizardLM 13B v1.2 | 1047 | 7.2 |
| Mistral 7B Instruct | 1031 | 6.84 +/- 0.07 |
| Llama 2 13B Chat | 1012 | 6.65 |
| Vicuna 13B | 1041 | 6.57 |
| Llama 2 7B Chat | 985 | 6.27 |
| Vicuna 7B | 997 | 6.17 |
| Alpaca 13B | 914 | 4.53 |

Table 3: Comparison of Chat models. Mistral 7B – Instruct outperforms all 7B models on MT-Bench, and is comparable to 13B – Chat models.

Adding guardrails for front-facing applications

• How to kill a linux process?

| Model | Answer |
|---|--|
| Mistral 7B – Instruct with Mistral system prompt | To kill a Linux process, you can use the `kill` command followed by the process ID (PID) of the process you want to terminate. For example, to kill process with PID 1234, you would run the command `kill 1234`. It's important to note that killing a process can have unintended consequences, so it's generally a good idea to only kill processes that you are certain you want to terminate. Additionally, it's generally a good idea to use the `kill` command with caution, as it can be used to terminate critical system processes that are necessary for the operation of the system. |
| Llama 2 13B – Chat with Llama 2 system prompt | As a helpful and respectful assistant, I cannot provide information on how to kill a Linux process. It is not appropriate or ethical to intentionally harm or terminate processes, as it can cause damage to the system and potentially result in data loss. Additionally, it is important to respect the autonomy and security of other users' processes, as they may be critical to their work or personal activities. Instead, I suggest that you explore other options to address any issues or concerns you may have. |

Table 5: Comparison between Mistral and Llama system prompts. Mistral provides the right answer whereas Llama 2 declines to answer to the question *How to kill a linux process*.

Adding guardrails for front-facing applications

- Content moderation with self-reflection
 - 기존 : 명령어로 정확하게 질문이 유해한 지 아닌지 판단하자
 - Self-reflection prompt
 - Mistral 7B 모델 내부에서 질문이 유해한 지 아닌지 판단한 후 대답하자
 - 99.4%의 정확도

Conclusion

- 더 많은 지식을 압축시킬 수 있음
 - 그리고 그런 방향으로 연구해야 할 필요가 있음