

EE508 Project - Phase_2

Asmita Mohanty

Shaun Almeida

Files need to be changed to disable KV Caching:

1. models.py
2. generation.py
3. inference.py

All changes to each file are highlighted in the red box.

Note: There are few masked comments in green in the code which are the debug prints done to understand the differences in the output results between with vs without KV Caching. They are not added to the final code. Debug print outputs shared at the end of this file.

1. models.py

Without KV Caching:

- Keys & Values of the *entire* token per sequence is computed
- Doesn't require to store the keys & values in cache anymore

With KV Caching:

- Keys & Values of the *current* token per sequence is only computed
- Stores the computed results in the cache & later assigns the entire cache that consists of the past tokens as well as the current tokens to the keys & values for downstream computation

```
if self.kv_caching:
    self.cache_k = self.cache_k.to(xq)
    self.cache_v = self.cache_v.to(xq)

    self.cache_k[:bsz, start_pos : start_pos + seqlen] = xk
    self.cache_v[:bsz, start_pos : start_pos + seqlen] = xv

    keys = self.cache_k[:bsz, : start_pos + seqlen]
    values = self.cache_v[:bsz, : start_pos + seqlen]

else:
    keys = xk
    values = xv
```

2. generation.py

- Tokens per sequence are passed as batched tensors

Without KV Caching:

- Pass the entire token per sequence to compute the logits

With KV Caching:

- Pass only the current token per sequence to compute the logits

```
if kv_caching:
    #for t,seq in enumerate(tokens[:, prev_pos:cur_pos]):
    #    print(f"W/ KV Cache, Tokens: {t}, seq:{seq.tolist()} of length:", len(seq.tolist()),
    #    # " decoded tokens:", tokenizer.decode(seq.tolist()))
    logits = self(tokens[:, prev_pos:cur_pos], prev_pos)
else:
    #for t,seq in enumerate(tokens[:, :cur_pos]):
    #    print(f"W/o KV Cache, Tokens: {t}, seq:{seq.tolist()} of length:", len(seq.tolist()),
    #    # " decoded tokens:", tokenizer.decode(seq.tolist()))
    logits = self(tokens[:, :cur_pos], prev_pos)
```

3. Inference.py

- Toggle the kv caching boolean flag

```
model.eval()
results = model.generate(tokenizer, prompts, max_gen_len=64, temperature=0.6, top_p=0.9, kv_caching=False, device=device)
```

In **benchmark_inference.py**: Toggle the batch_size & kv_caching flag to compare the KV caching results with different batch sizes.

```
if __name__ == "__main__":
    benchmark_inference(batch_size=1, input_len=256, output_len=64, kv_caching=False)
```

Output Results:

All outputs tested with input length = 256 tokens & output length = 32 tokens

KV Caching	Batch size = 1	Batch size = 8	Batch size = 16
With KV Cache	Peak Memory: 3071.57MB Run Time: 1.46s	Peak Memory: 4495.47MB Run Time: 2.02s	Peak Memory: 6133.8MB Run Time: 8.14s
Without KV Cache	Peak Memory: 3230.12MB Run Time: 4.15s	Peak Memory: 5755MB Run Time: 28.29s	Peak Memory: 8641.54MB Run Time: 52.3s

Output Logs:

Batch Size	With KV Cache	Without KV Cache
16	<pre> Reloaded tiktoken model from /home1/asmitamo/.llama/checkpoints/Llama3.2-1B/tokenizer.model #words: 128256 - BOS ID: 128000 - EOS ID: 128001 Batch size: 16 Input length: 256 tokens Output length: 32 tokens Inference time: 8.14 seconds Tokens per second: 62.88 Model weights memory usage: 2858.13 MB Peak memory usage: 6133.80 MB === Sample Output === time in a galaxy far away Once upon a time in a galaxy far away Once upon a time in a galaxy </pre>	<pre> Reloaded tiktoken model from /home1/asmitamo/.llama/checkpoints/Llama3.2-1B/tokenizer.model #words: 128256 - BOS ID: 128000 - EOS ID: 128001 Batch size: 16 Input length: 256 tokens Output length: 32 tokens Inference time: 52.30 seconds Tokens per second: 9.79 Model weights memory usage: 2858.13 MB Peak memory usage: 8641.54 MB === Sample Output === time in a galaxy far away Once upon a time in a galaxy far away Once upon a time in a galaxy </pre>
8	<pre> Reloaded tiktoken model from /home1/asmitamo/.llama/checkpoints/Llama3.2-1B/tokenizer.model #words: 128256 - BOS ID: 128000 - EOS ID: 128001 Batch size: 8 Input length: 256 tokens Output length: 32 tokens Inference time: 2.02 seconds Tokens per second: 126.76 Model weights memory usage: 2858.13 MB Peak memory usage: 4495.47 MB === Sample Output === time in a galaxy far away Once upon a time in a galaxy far away Once upon a time in a galaxy far away Once upon a </pre>	<pre> Reloaded tiktoken model from /home1/asmitamo/.llama/checkpoints/Llama3.2-1B/tokenizer.model #words: 128256 - BOS ID: 128000 - EOS ID: 128001 Batch size: 8 Input length: 256 tokens Output length: 32 tokens Inference time: 28.29 seconds Tokens per second: 9.05 Model weights memory usage: 2858.13 MB Peak memory usage: 5755.00 MB === Sample Output === time in a galaxy far away Once upon a time in a galaxy far away Once upon a time in a galaxy </pre>
1	<pre> Reloaded tiktoken model from /home1/asmitamo/.llama/checkpoints/Llama3.2-1B/tokenizer.model #words: 128256 - BOS ID: 128000 - EOS ID: 128001 Batch size: 1 Input length: 256 tokens Output length: 32 tokens Inference time: 1.46 seconds Tokens per second: 21.91 Model weights memory usage: 2858.13 MB Peak memory usage: 3071.57 MB === Sample Output === time in a galaxy far away Once upon a time in a galaxy far away Once upon a time in a galaxy </pre>	<pre> Reloaded tiktoken model from /home1/asmitamo/.llama/checkpoints/Llama3.2-1B/tokenizer.model #words: 128256 - BOS ID: 128000 - EOS ID: 128001 Batch size: 1 Input length: 256 tokens Output length: 32 tokens Inference time: 4.15 seconds Tokens per second: 7.72 Model weights memory usage: 2858.13 MB Peak memory usage: 3230.12 MB === Sample Output === time in a galaxy far away Once upon a time in a galaxy far away Once upon a time in a galaxy </pre>

Additional Data

Debug Outputs explaining with/without KV Caching differences

1. Given prompts (same for both w/ KV vs w/o KV cache):

```

Prompts: ['I believe the meaning of life is', 'Simply put, the theory of relativity states that ', 'A brief message congratulating the team on the launch']
Prompt tokens: [[128000, 40, 4510, 279, 7438, 315, 2324, 374], [128000, 61346, 2231, 11, 279, 10334, 315, 1375, 44515, 5415, 128001], [128000, 128001]]
length of token: 8 decoded tokens: <|begin_of_text|>I believe the meaning of life is
length of token: 12 decoded tokens: <|begin_of_text|>Simply put, the theory of relativity states that
length of token: 21 decoded tokens: <|begin_of_text|>A brief message congratulating the team on the launch:

Hi everyone,

I just
length of token: 41 decoded tokens: <|begin_of_text|>Translate English to French:

sea otter => loutre de mer
peppermint => menthe poivrée
plush giraffe => girafe peluche
cheese =>

```

- During token generation, we iterate over the batched tensor which contains sequences of varying lengths. Starting the 2nd iteration, inference “with KV cache” will generate only 1 token while inference “without KV cache” will output the newly generated token along with the previous/past tokens. Debug outputs shown for first 3 iterations.

Legend:

Pink - Past/Previous tokens

Red - Newly generated tokens

Iteration #	With KV Cache	Without KV Cache
1	<pre>cur_pos: 0 prev_pos: 0 w/ KV Cache, Tokens: 0, seq:[128000, 40, 4510, 279, 7438, 315, 2324, 374] of length: 8 decoded tokens: <[begin_of_text]>I believe the meaning of life is w/ KV Cache, Tokens: 1, seq:[128000, 61346, 2231, 11, 279, 10334, 315, 1375] of length: 8 decoded tokens: <[begin_of_text]>Simply put, the theory of rel w/ KV Cache, Tokens: 2, seq:[128000, 32, 10015, 1984, 40588, 15853, 279, 2128] of length: 8 decoded tokens: <[begin_of_text]>A brief message congratulating the team w/ KV Cache, Tokens: 3, seq:[128000, 28573, 6498, 311, 8753, 512, 1827, 286] of length: 8 decoded tokens: <[begin_of_text]>Translate English to French:</pre>	<pre>cur_pos: 0 prev_pos: 0 w/o KV Cache, Tokens: 0, seq:[128000, 40, 4510, 279, 7438, 315, 2324, 374] of length: 8 decoded tokens: <[begin_of_text]>I believe the meaning of life is w/o KV Cache, Tokens: 1, seq:[128000, 61346, 2231, 11, 279, 10334, 315, 1375] of length: 8 decoded tokens: <[begin_of_text]>Simply put, the theory of rel w/o KV Cache, Tokens: 2, seq:[128000, 32, 10015, 1984, 40588, 15853, 279, 2128] of length: 8 decoded tokens: <[begin_of_text]>A brief message congratulating the team w/o KV Cache, Tokens: 3, seq:[128000, 28573, 6498, 311, 8753, 512, 1827, 286] of length: 8 decoded tokens: <[begin_of_text]>Translate English to French:</pre>
2	<pre>cur_pos: 9 prev_pos: 8 w/ KV Cache, Tokens: 0, seq:[311] of length: 1 decoded tokens: to w/ KV Cache, Tokens: 1, seq:[44515] of length: 1 decoded tokens: activity w/ KV Cache, Tokens: 2, seq:[389] of length: 1 decoded tokens: on w/ KV Cache, Tokens: 3, seq:[9581] of length: 1 decoded tokens: sea</pre>	<pre>cur_pos: 9 prev_pos: 8 w/o KV Cache, Tokens: 0, seq:[128000, 40, 4510, 279, 7438, 315, 2324, 374, 311] of length: 9 decoded tokens: [begin_of_text]I believe the meaning of life is to w/o KV Cache, Tokens: 1, seq:[128000, 61346, 2231, 11, 279, 10334, 315, 1375, 44515] of length: 9 decoded tokens: [begin_of_text]Simply put, the theory of relativity on w/o KV Cache, Tokens: 2, seq:[128000, 32, 10015, 1984, 40588, 15853, 279, 2128, 389] of length: 9 decoded tokens: [begin_of_text]A brief message congratulating the team on w/o KV Cache, Tokens: 3, seq:[128000, 28573, 6498, 311, 8753, 512, 1827, 286, 9581] of length: 9 decoded tokens: [begin_of_text]Translate English to French: sea</pre>
3	<pre>cur_pos: 10 prev_pos: 9 w/ KV Cache, Tokens: 0, seq:[3974] of length: 1 decoded tokens: live w/ KV Cache, Tokens: 1, seq:[5415] of length: 1 decoded tokens: states w/ KV Cache, Tokens: 2, seq:[279] of length: 1 decoded tokens: the w/ KV Cache, Tokens: 3, seq:[14479] of length: 1 decoded tokens: ot</pre>	<pre>cur_pos: 10 prev_pos: 9 w/o KV Cache, Tokens: 0, seq:[128000, 40, 4510, 279, 7438, 315, 2324, 374, 311, 3974] of length: 10 decoded tokens: [begin_of_text]I believe the meaning of life is to live w/o KV Cache, Tokens: 1, seq:[128000, 61346, 2231, 11, 279, 10334, 315, 1375, 44515, 5415] of length: 10 decoded tokens: [begin_of_text]Simply put, the theory of relativity states w/o KV Cache, Tokens: 2, seq:[128000, 32, 10015, 1984, 40588, 15853, 279, 2128, 389, 279] of length: 10 decoded tokens: [begin_of_text]A brief message congratulating the team the w/o KV Cache, Tokens: 3, seq:[128000, 28573, 6498, 311, 8753, 512, 1827, 286, 9581, 14479] of length: 10 decoded tokens: [begin_of_text]Translate English to French: sea ot</pre>