

EE 599: Systems for Machine Learning
Final Year Project - Phase 2

Efficient Fine-tuning of LLM on a Single GPU

Sheik Dhawood Ashfaq Asick Ali
Sam Devavaram Jebaraj

Time taken with KV caching: 48.36289 seconds

Time taken without KV caching: 93.8209 seconds

The presence of KV caching likely speeds up inference by reducing the time required for key-value pair retrieval, enhancing overall efficiency. The variation in outputs could be attributed to differences in GPU computation kernels or Python package versions. These variances can lead to differences in floating-point calculations, potentially affecting the output.

Given Prompts =

```
[
    # For these prompts, the expected answer is the natural
continuation of the prompt
    "I believe the meaning of life is",
    "Simply put, the theory of relativity states that ",
    """"A brief message congratulating the team on the launch:

Hi everyone,

I just """,
    # Few shot prompt (providing a few examples before asking
model to complete more);
    """"Translate English to French:

sea otter => loutre de mer
peppermint => menthe poivrée
plush girafe => girafe peluche
cheese =>""",
]
```

Changes in code:

In Model Py: Commenting off the KV Caching in the Attention Module and directly compute keys and values from the input tensors.

```

188     super().__init__()
189     self.n_kv_heads = args.n_heads if args.n_kv_heads is None else args.n_kv_heads
190     self.n_local_heads = args.n_heads
191     self.n_local_kv_heads = self.n_kv_heads
192     self.n_rep = self.n_local_heads // self.n_local_kv_heads
193     self.head_dim = args.dim // args.n_heads
194
195     self.wq = nn.Linear(args.dim, args.n_heads * self.head_dim, bias=False)
196     self.wk = nn.Linear(args.dim, self.n_kv_heads * self.head_dim, bias=False)
197     self.wv = nn.Linear(args.dim, self.n_kv_heads * self.head_dim, bias=False)
198     self.wo = nn.Linear(args.n_heads * self.head_dim, args.dim, bias=False)
199
200     '''self.cache_k = torch.zeros(
201         (
202             args.max_batch_size,
203             args.max_seq_len,
204             self.n_local_kv_heads,
205             self.head_dim,
206         )
207     ).cuda()
208     self.cache_v = torch.zeros(
209         (
210             args.max_batch_size,
211             args.max_seq_len,
212             self.n_local_kv_heads,
213             self.head_dim,
214         )
215     ).cuda()'''
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246     '''self.cache_k = self.cache_k.to(xq)
247     self.cache_v = self.cache_v.to(xq)
248
249     self.cache_k[:bsz, start_pos : start_pos + seqlen] = xk
250     self.cache_v[:bsz, start_pos : start_pos + seqlen] = xv
251
252     keys = self.cache_k[:bsz, : start_pos + seqlen]
253     values = self.cache_v[:bsz, : start_pos + seqlen]'''
254     keys = xk
255     values = xv

```

In Generation Py: Making the previous position to be 0, thus it doesn't update and no slicing occurs.

```

22         for cur_pos in range(min_prompt_len, total_len):
23             with torch.no_grad():
24                 logits = self(tokens[:, prev_pos:cur_pos], prev_pos)
25                 if temperature > 0:
26                     probs = torch.softmax(logits[:, -1] / temperature, dim=-1)
27                     next_token = sample_top_p(probs, top_p)
28                 else:
29                     next_token = torch.argmax(logits[:, -1], dim=-1)
30
31             next_token = next_token.reshape(-1)
32             # only replace token if prompt has already been generated
33             next_token = torch.where(
34                 input_text_mask[:, cur_pos], tokens[:, cur_pos], next_token
35             )
36             tokens[:, cur_pos] = next_token
37
38             eos_reached |= (~input_text_mask[:, cur_pos]) & (
39                 next_token == tokenizer.eos_id
40             )
41
42         #prev_pos = cur_pos

```

Without KV Caching	With KV Caching
<p>I believe the meaning of life is > to learn to love. Love is not a feeling. It is a decision. I believe the meaning of life is to learn to love. Love is not a feeling. It is a decision. It is a commitment. It is a conscious choice of the will and the intellect. There are many</p> <p>=====</p> <p>Simply put, the theory of relativity states that > 1) the speed of light is constant for all observers and 2) the laws of physics are the same for all observers. The theory of relativity is a very important concept in physics, but it is also one of the most misunderstood. There are a lot of misconceptions about</p> <p>=====</p> <p>A brief message congratulating the team on the launch:</p> <p>Hi everyone,</p> <p>I just</p> <p>></p> <p>Google your website.</p> <p>I hope you enjoy the new look and feel.</p> <p>I'll be in touch soon to discuss your next project.</p> <p>Best</p> <p>=====</p> <p>Translate English to French:</p>	<p>I believe the meaning of life is > to learn to love. Love is not a feeling. It is a decision. I believe the meaning of life is to learn to love. Love is not a feeling. It is a decision. It is a commitment. It is a conscious choice of the will and the intellect. There are many</p> <p>=====</p> <p>Simply put, the theory of relativity states that > 1) the speed of light is constant for all observers and 2) the laws of physics are the same for all observers. The theory of relativity is a very important concept in physics, but it is also one of the most misunderstood. There are a lot of misconceptions about</p> <p>=====</p> <p>A brief message congratulating the team on the launch:</p> <p>Hi everyone,</p> <p>I just</p> <p>></p> <p>Google your website.</p> <p>I hope you enjoy the new look and feel.</p> <p>I'll be in touch soon to discuss your next project.</p> <p>Best</p> <p>=====</p> <p>Translate English to French:</p>

<pre> sea otter => loutre de mer peppermint => menthe poivrée plush girafe => girafe peluche cheese => > fromage penguin => pinguin handbag => sac à main mug => tasse chocolate => chocolat chocolate => chocolat chocolate => chocolat chocolate => choc ===== </pre>	<pre> sea otter => loutre de mer peppermint => menthe poivrée plush girafe => girafe peluche cheese => > fromage penguin => pinguin handbag => sac à main mug => tasse toothpaste => dentifrice t-shirt => tee-shirt pencil => crayon parrot => perroquet ===== </pre>
-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Custom prompt: “Best Sunset hike in Los Angeles”

Time taken with KV caching: 32.03430461883545 secs

Time taken without KV caching: 35.742812633514404 secs

Both Outputs:

Best Hiking place in Los Angeles for sunset

>

L.A. is known for its beautiful sunsets, and there are many hiking trails in the city that offer stunning views of the setting sun. Here are some of the best hiking places in Los Angeles for sunset:

1. Griffith Park: Griffith Park is a