

EE599 - Systems for ML Phase 2

Team: Prasanna Akolkar, Sreetama Sarkar

Details of the changes:

File name - model.py

1 - class Attention(nn.Module):

Code to be removed:

```
self.cache_k = torch.zeros(
    (
        args.max_batch_size,
        args.max_seq_len,
        self.n_local_kv_heads,
        self.head_dim,
    )
).cuda()
self.cache_v = torch.zeros(
    (
        args.max_batch_size,
        args.max_seq_len,
        self.n_local_kv_heads,
        self.head_dim,
    )
).cuda()
```

forward method of attention class

- removed the start_pos from the parameters
- removed the following code

```
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]

repeat k/v heads if n_kv_heads < n_heads
keys = repeat_kv(keys, self.n_rep) # (bs, cache_len + seqlen, n_local_heads,
head_dim)
values = repeat_kv(values, self.n_rep) # (bs, cache_len + seqlen,
n_local_heads, head_dim)

keys = keys.transpose(1, 2) # (bs, n_local_heads, cache_len + seqlen, head_dim)
values = values.transpose(1, 2) # (bs, n_local_heads, cache_len + seqlen,
head_dim)

scores = torch.matmul(xq, keys.transpose(2, 3)) / math.sqrt(self.head_dim)

```

- updated the following code using xk, xq, xv

```

xq = xq.transpose(1, 2) # (bs, n_local_heads, seqlen, head_dim)
xk = xk.transpose(1, 2) # (bs, n_kv_heads, seqlen, head_dim) new
xv = xv.transpose(1, 2) # (bs, n_kv_heads, seqlen, head_dim) new

scores = torch.matmul(xq, xk.transpose(2, 3)) / math.sqrt(self.head_dim)
if mask is not None:
    scores = scores + mask # (bs, n_local_heads, seqlen, cache_len + seqlen)
scores = F.softmax(scores.float(), dim=-1).type_as(xq)
output = torch.matmul(scores, xv) # (bs, n_local_heads, seqlen, head_dim)
output = output.transpose(1, 2).contiguous().view(bsz, seqlen, -1)

```

2 - class TransformerBlock(nn.Module):

updated the forward method of transformer block

- removed start_pos from function parameters
- updated the parameters while making call to the forward method of attention block (without start_pos)

```
h = x + self.attention.forward(
    self.attention_norm(x), freqs_cis, mask
)
```

3 - class Llama(Generation):

updated the forward method of the LLama block

- removed start_pos from the function parameters
- updated freq_cis without start_pos

```
freqs_cis = self.freqs_cis[:seqlen]
```

```
for layer in self.layers:
    h = layer(h, freqs_cis, mask)
```

updated above code without start_pos

Apart from these, we can remove the repeat_kv function since it is no longer used.

These are the changes in the model.py

File name - generation.py

We make only one line of change in this

originally:

```
logits = self(tokens[:, prev_pos:cur_pos], prev_pos)
```

updated:

```
logits = self(tokens[:, :cur_pos])
```

Apart from these two files, no other changes are required.

The code is running without bugs and generated outputs according to the original model.

Following are the prompts (no changes made to the ones provided)

```
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 poivr e  
    plush girafe => girafe peluche  
    cheese =>""",  
]
```

Original Code - With KV Caching Output

```

Singularity> python inference.py
/home1/pakolkar/.local/lib/python3.9/site-packages/torch/__init__.py:696: UserWarning: torch.set_default_tensor_type() is deprecated as
of PyTorch 2.1, please use torch.set_default_dtype() and torch.set_default_device() as alternatives. (Triggered internally at ../torch/c
src/tensor/python_tensor.cpp:451.)
  C._set_default_tensor_type(t)
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

=====

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

=====

A brief message congratulating the team on the launch:

    Hi everyone,

    I just

>

    <a href="https://www.google.com">Google</a>

    your website.

    I hope you enjoy the new look and feel.

    I'll be in touch soon to discuss your next project.

    Best

=====

```

Translate English to French:

```

    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

```

Updated Code - Without KV Caching Output

```

Singularity> python inference.py
/home1/pakolkar/.local/lib/python3.9/site-packages/torch/_init_.py:696: UserWarning: torch.set_default_tensor_type() is deprecated as
of PyTorch 2.1, please use torch.set_default_dtype() and torch.set_default_device() as alternatives. (Triggered internally at ../torch/c
src/tensor/python_tensor.cpp:451.)
  C._set_default_tensor_type(t)
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

=====

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

=====

A brief message congratulating the team on the launch:

    Hi everyone,

    I just

>

    <a href="https://www.google.com">Google</a>

    your website.

    I hope you enjoy the new look and feel.

    I'll be in touch soon to discuss your next project.

```

```

    Best

=====

Translate English to French:

    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

```