

"""
Full definition of a GPT Language Model, all of it in this single file.

References:

1) the official GPT-2 TensorFlow implementation released by OpenAI:

<https://github.com/openai/gpt-2/blob/master/src/model.py>

2) huggingface/transformers PyTorch implementation:

https://github.com/huggingface/transformers/blob/main/src/transformers/models/gpt2/modeling_gpt2.py

"""

```
import math
import inspect
from dataclasses import dataclass

import torch
import torch.nn as nn
from torch.nn import functional as F
```

```
class LayerNorm(nn.Module):
```

""" LayerNorm but with an optional bias. PyTorch doesn't support simply bias=False """

```
def __init__(self, ndim, bias):
    super().__init__()
    self.weight = nn.Parameter(torch.ones(ndim))
    self.bias = nn.Parameter(torch.zeros(ndim)) if bias else None
```

'''

In the computational graph of layer norm, it just applies layer normalization to input.

'''

```
def forward(self, input):
    return F.layer_norm(input, self.weight.shape, self.weight, self.bias, 1e-5)
```

```
class CausalSelfAttention(nn.Module):
```

```
def __init__(self, config):
    super().__init__()
    assert config.n_embd % config.n_head == 0
    # key, query, value projections for all heads, but in a batch
    self.c_attn = nn.Linear(config.n_embd, 3 * config.n_embd, bias=config.bias)
    # output projection
    self.c_proj = nn.Linear(config.n_embd, config.n_embd, bias=config.bias)
    # regularization
    self.attn_dropout = nn.Dropout(config.dropout)
    self.resid_dropout = nn.Dropout(config.dropout)
    self.n_head = config.n_head
    self.n_embd = config.n_embd
    self.dropout = config.dropout
    # flash attention make GPU go brrrrr but support is only in PyTorch >= 2.0
    self.flash = hasattr(torch.nn.functional, 'scaled_dot_product_attention')
    if not self.flash:
        print("WARNING: using slow attention. Flash Attention requires PyTorch >= 2.0")
        # causal mask to ensure that attention is only applied to the left in the input
        sequence
        self.register_buffer("bias", torch.tril(torch.ones(config.block_size,
config.block_size))
                                .view(1, 1, config.block_size, config.block_size))
```

'''

The CausalSelfAttention class implements the self-attention mechanism, which includes matrix multiplications and masking operations that are part of the computation graph. Here, matrix multiplications (@ operator), attention weight calculation (softmax), and dropout (self.resid_dropout) all contribute to the computation graph.

Below is the overall flow of the computational graph for CasualSelfAttention:

- The input is split into queries, keys, and values, which are then used in the attention mechanism.
- Depending on whether Flash Attention is available, either a fast or a manual implementation of attention is used.

- The resulting tensor is processed by dropout and an output projection layer to complete the self-attention step.

```
'''
def forward(self, x):
    B, T, C = x.size() # batch size, sequence length, embedding dimensionality (n_embd)

    # calculate query, key, values for all heads in batch and move head forward to be the
    batch dim
    q, k, v = self.c_attn(x).split(self.n_embd, dim=2)
    k = k.view(B, T, self.n_head, C // self.n_head).transpose(1, 2) # (B, nh, T, hs)
    q = q.view(B, T, self.n_head, C // self.n_head).transpose(1, 2) # (B, nh, T, hs)
    v = v.view(B, T, self.n_head, C // self.n_head).transpose(1, 2) # (B, nh, T, hs)

    # causal self-attention; Self-attend: (B, nh, T, hs) x (B, nh, hs, T) -> (B, nh, T, T)
    if self.flash:
        # efficient attention using Flash Attention CUDA kernels
        y = torch.nn.functional.scaled_dot_product_attention(q, k, v, attn_mask=None,
dropout_p=self.dropout if self.training else 0, is_causal=True)
    else:
        # manual implementation of attention
        att = (q @ k.transpose(-2, -1)) * (1.0 / math.sqrt(k.size(-1)))
        att = att.masked_fill(self.bias[:, :, :T, :T] == 0, float('-inf'))
        att = F.softmax(att, dim=-1)
        att = self.attn_dropout(att)
        y = att @ v # (B, nh, T, T) x (B, nh, T, hs) -> (B, nh, T, hs)
    y = y.transpose(1, 2).contiguous().view(B, T, C) # re-assemble all head outputs side by
    side

    # output projection
    y = self.resid_dropout(self.c_proj(y))
    return y
```

```
class MLP(nn.Module):
```

```
def __init__(self, config):
    super().__init__()
    self.c_fc = nn.Linear(config.n_embd, 4 * config.n_embd, bias=config.bias)
    self.gelu = nn.GELU()
    self.c_proj = nn.Linear(4 * config.n_embd, config.n_embd, bias=config.bias)
    self.dropout = nn.Dropout(config.dropout)
```

'''

The MLP class defines a multi-layer perceptron, which includes linear transformations, GELU activation, and dropout.

These operations are also part of below the computation graph. The linear transformations (self.c_fc, self.c_proj), non-linear activation (self.gelu), and dropout (self.dropout) form part of the forward computation.

```
'''
def forward(self, x):
    x = self.c_fc(x)
    x = self.gelu(x)
    x = self.c_proj(x)
    x = self.dropout(x)
    return x
```

```
class Block(nn.Module):
```

```
def __init__(self, config):
    super().__init__()
    self.ln_1 = LayerNorm(config.n_embd, bias=config.bias)
    self.attn = CausalSelfAttention(config)
    self.ln_2 = LayerNorm(config.n_embd, bias=config.bias)
    self.mlp = MLP(config)
```

'''

In the Block class, the computation graph is built by applying the LayerNorm, CausalSelfAttention, and MLP operations to the input tensor dL. This is a critical part of the transformer block's

forward pass.

The operations applied hereâ€self.attn (which calls the attention mechanism), self.mlp (the multi-layer perceptron), and self.ln_1 and self.ln_2 (layer normalization)â€are part of the computation graph.

'''

```
def forward(self, x):
    x = x + self.attn(self.ln_1(x))
    x = x + self.mlp(self.ln_2(x))
    return x
```

@dataclass

class GPTConfig:

```
    block_size: int = 1024
    vocab_size: int = 50304 # GPT-2 vocab_size of 50257, padded up to nearest multiple of 64 for
efficiency
    n_layer: int = 12
    n_head: int = 12
    n_embd: int = 768
    dropout: float = 0.0
    bias: bool = True # True: bias in Linears and LayerNorms, like GPT-2. False: a bit better and
faster
```

class GPT(nn.Module):

```
    def __init__(self, config):
        super().__init__()
        assert config.vocab_size is not None
        assert config.block_size is not None
        self.config = config

        self.transformer = nn.ModuleDict(dict(
            wte = nn.Embedding(config.vocab_size, config.n_embd),
            wpe = nn.Embedding(config.block_size, config.n_embd),
            drop = nn.Dropout(config.dropout),
            h = nn.ModuleList([Block(config) for _ in range(config.n_layer)]),
            ln_f = LayerNorm(config.n_embd, bias=config.bias),
        ))
        self.lm_head = nn.Linear(config.n_embd, config.vocab_size, bias=False)
        # with weight tying when using torch.compile() some warnings get generated:
        # "UserWarning: functional_call was passed multiple values for tied weights.
        # This behavior is deprecated and will be an error in future versions"
        # not 100% sure what this is, so far seems to be harmless. TODO investigate
        self.transformer.wte.weight = self.lm_head.weight #
https://paperswithcode.com/method/weight-tying

        # init all weights
        self.apply(self._init_weights)
        # apply special scaled init to the residual projections, per GPT-2 paper
        for pn, p in self.named_parameters():
            if pn.endswith('c_proj.weight'):
                torch.nn.init.normal_(p, mean=0.0, std=0.02/math.sqrt(2 * config.n_layer))

        # report number of parameters
        print("number of parameters: %.2fM" % (self.get_num_params()/1e6,))

    def get_num_params(self, non_embedding=True):
        """
        Return the number of parameters in the model.
        For non-embedding count (default), the position embeddings get subtracted.
        The token embeddings would too, except due to the parameter sharing these
        params are actually used as weights in the final layer, so we include them.
        """
        n_params = sum(p.numel() for p in self.parameters())
        if non_embedding:
            n_params -= self.transformer.wpe.weight.numel()
        return n_params

    def _init_weights(self, module):
```

```

if isinstance(module, nn.Linear):
    torch.nn.init.normal_(module.weight, mean=0.0, std=0.02)
    if module.bias is not None:
        torch.nn.init.zeros_(module.bias)
elif isinstance(module, nn.Embedding):
    torch.nn.init.normal_(module.weight, mean=0.0, std=0.02)

```

'''
The primary computation graph for the entire model is built within the forward method of the GPT class.

This method handles token embeddings, position embeddings, and the sequential application of transformer blocks.

- The forward pass starts by embedding the input idx using the token embeddings (self.transformer.wte) and position embeddings (self.transformer.wpe).*
- These embeddings are passed through multiple blocks of the model (defined in Block).*
- Finally, the output is passed through a linear layer (self.lm_head) to produce the logits.*
- If targets are provided, the cross-entropy loss is computed.*
- This method defines the high-level structure of the computation graph for a forward pass of the GPT model.*

```

'''
def forward(self, idx, targets=None):
    device = idx.device
    b, t = idx.size()
    assert t <= self.config.block_size, f"Cannot forward sequence of length {t}, block size is only {self.config.block_size}"
    pos = torch.arange(0, t, dtype=torch.long, device=device) # shape (t)

    # forward the GPT model itself
    tok_emb = self.transformer.wte(idx) # token embeddings of shape (b, t, n_embd)
    pos_emb = self.transformer.wpe(pos) # position embeddings of shape (t, n_embd)
    x = self.transformer.drop(tok_emb + pos_emb)
    for block in self.transformer.h:
        x = block(x)
    x = self.transformer.ln_f(x)

    if targets is not None:
        # if we are given some desired targets also calculate the loss
        logits = self.lm_head(x)
        loss = F.cross_entropy(logits.view(-1, logits.size(-1)), targets.view(-1),
ignore_index=-1)
    else:
        # inference-time mini-optimization: only forward the lm_head on the very last
position
        logits = self.lm_head(x[:, [-1], :]) # note: using list [-1] to preserve the time dim
        loss = None

    return logits, loss

```

```

def crop_block_size(self, block_size):
    # model surgery to decrease the block size if necessary
    # e.g. we may load the GPT2 pretrained model checkpoint (block size 1024)
    # but want to use a smaller block size for some smaller, simpler model
    assert block_size <= self.config.block_size
    self.config.block_size = block_size
    self.transformer.wpe.weight = nn.Parameter(self.transformer.wpe.weight[:block_size])
    for block in self.transformer.h:
        if hasattr(block.attn, 'bias'):
            block.attn.bias = block.attn.bias[:, :, :block_size, :block_size]

```

```

@classmethod
def from_pretrained(cls, model_type, override_args=None):
    assert model_type in {'gpt2', 'gpt2-medium', 'gpt2-large', 'gpt2-xl'}
    override_args = override_args or {} # default to empty dict
    # only dropout can be overridden see more notes below
    assert all(k == 'dropout' for k in override_args)
    from transformers import GPT2LMHeadModel
    print("loading weights from pretrained gpt: %s" % model_type)

```

```

# n_layer, n_head and n_embd are determined from model_type
config_args = {
    'gpt2': dict(n_layer=12, n_head=12, n_embd=768), # 124M params
    'gpt2-medium': dict(n_layer=24, n_head=16, n_embd=1024), # 350M params
    'gpt2-large': dict(n_layer=36, n_head=20, n_embd=1280), # 774M params
    'gpt2-xl': dict(n_layer=48, n_head=25, n_embd=1600), # 1558M params
} [model_type]
print("forcing vocab_size=50257, block_size=1024, bias=True")
config_args['vocab_size'] = 50257 # always 50257 for GPT model checkpoints
config_args['block_size'] = 1024 # always 1024 for GPT model checkpoints
config_args['bias'] = True # always True for GPT model checkpoints
# we can override the dropout rate, if desired
if 'dropout' in override_args:
    print(f"overriding dropout rate to {override_args['dropout']}")
    config_args['dropout'] = override_args['dropout']
# create a from-scratch initialized minGPT model
config = GPTConfig(**config_args)
model = GPT(config)
sd = model.state_dict()
sd_keys = sd.keys()
sd_keys = [k for k in sd_keys if not k.endswith('.attn.bias')] # discard this mask /
buffer, not a param

# init a huggingface/transformers model
model_hf = GPT2LMHeadModel.from_pretrained(model_type)
sd_hf = model_hf.state_dict()

# copy while ensuring all of the parameters are aligned and match in names and shapes
sd_keys_hf = sd_hf.keys()
sd_keys_hf = [k for k in sd_keys_hf if not k.endswith('.attn.masked_bias')] # ignore
these, just a buffer
sd_keys_hf = [k for k in sd_keys_hf if not k.endswith('.attn.bias')] # same, just the
mask (buffer)
transposed = ['attn.c_attn.weight', 'attn.c_proj.weight', 'mlp.c_fc.weight',
'mlp.c_proj.weight']
# basically the openai checkpoints use a "Conv1D" module, but we only want to use a
vanilla Linear
# this means that we have to transpose these weights when we import them
assert len(sd_keys_hf) == len(sd_keys), f"mismatched keys: {len(sd_keys_hf)} !=
{len(sd_keys)}"
for k in sd_keys_hf:
    if any(k.endswith(w) for w in transposed):
        # special treatment for the Conv1D weights we need to transpose
        assert sd_hf[k].shape[:-1] == sd[k].shape
        with torch.no_grad():
            sd[k].copy_(sd_hf[k].t())
    else:
        # vanilla copy over the other parameters
        assert sd_hf[k].shape == sd[k].shape
        with torch.no_grad():
            sd[k].copy_(sd_hf[k])

return model

def configure_optimizers(self, weight_decay, learning_rate, betas, device_type):
    # start with all of the candidate parameters
    param_dict = {pn: p for pn, p in self.named_parameters()}
    # filter out those that do not require grad
    param_dict = {pn: p for pn, p in param_dict.items() if p.requires_grad}
    # create optim groups. Any parameters that is 2D will be weight decayed, otherwise no.
    # i.e. all weight tensors in matmuls + embeddings decay, all biases and layernorms don't.
    decay_params = [p for n, p in param_dict.items() if p.dim() >= 2]
    nodecay_params = [p for n, p in param_dict.items() if p.dim() < 2]
    optim_groups = [
        {'params': decay_params, 'weight_decay': weight_decay},
        {'params': nodecay_params, 'weight_decay': 0.0}
    ]
    num_decay_params = sum(p.numel() for p in decay_params)

```

```

num_nodecay_params = sum(p.numel() for p in nodecay_params)
print(f"num decayed parameter tensors: {len(decay_params)}, with {num_decay_params:,}
parameters")
print(f"num non-decayed parameter tensors: {len(nodecay_params)}, with
{num_nodecay_params:,} parameters")
# Create AdamW optimizer and use the fused version if it is available
fused_available = 'fused' in inspect.signature(torch.optim.AdamW).parameters
use_fused = fused_available and device_type == 'cuda'
extra_args = dict(fused=True) if use_fused else dict()
optimizer = torch.optim.AdamW(optimizer_groups, lr=learning_rate, betas=betas, **extra_args)
print(f"using fused AdamW: {use_fused}")

return optimizer

def estimate_mfu(self, fwdbwd_per_iter, dt):
    """ estimate model flops utilization (MFU) in units of A100 bfloat16 peak FLOPS """
    # first estimate the number of flops we do per iteration.
    # see PaLM paper Appendix B as ref: https://arxiv.org/abs/2204.02311
    N = self.get_num_params()
    cfg = self.config
    L, H, Q, T = cfg.n_layer, cfg.n_head, cfg.n_embd//cfg.n_head, cfg.block_size
    flops_per_token = 6*N + 12*L*H*Q*T
    flops_per_fwdbwd = flops_per_token * T
    flops_per_iter = flops_per_fwdbwd * fwdbwd_per_iter
    # express our flops throughput as ratio of A100 bfloat16 peak flops
    flops_achieved = flops_per_iter * (1.0/dt) # per second
    flops_promised = 312e12 # A100 GPU bfloat16 peak flops is 312 TFLOPS
    mfu = flops_achieved / flops_promised
    return mfu

@torch.no_grad()
def generate(self, idx, max_new_tokens, temperature=1.0, top_k=None):
    """
    Take a conditioning sequence of indices idx (LongTensor of shape (b,t)) and complete
    the sequence max_new_tokens times, feeding the predictions back into the model each time.
    Most likely you'll want to make sure to be in model.eval() mode of operation for this.
    """
    for _ in range(max_new_tokens):
        # if the sequence context is growing too long we must crop it at block_size
        idx_cond = idx if idx.size(1) <= self.config.block_size else idx[:,-self.config.block_size:]
        # forward the model to get the logits for the index in the sequence
        logits, _ = self(idx_cond)
        # pluck the logits at the final step and scale by desired temperature
        logits = logits[:, -1, :] / temperature
        # optionally crop the logits to only the top k options
        if top_k is not None:
            v, _ = torch.topk(logits, min(top_k, logits.size(-1)))
            logits[logits < v[:, [-1]]] = -float('Inf')
        # apply softmax to convert logits to (normalized) probabilities
        probs = F.softmax(logits, dim=-1)
        # sample from the distribution
        idx_next = torch.multinomial(probs, num_samples=1)
        # append sampled index to the running sequence and continue
        idx = torch.cat((idx, idx_next), dim=1)

    return idx

```