

batch = B
seq_len = L
dim = D
n_heads = Nh
n_kv_heads = Nkv
hidden_dim = Dh
n_layer = Nl
vocab_size = V

Llama model architecture accoding to facebookresearch/llama/llama_model.py and karpathy/llama2.c/model.py

