
Karpathy Inference for Llama-2

Ramesh Subramonian

December 28, 2025

Abstract

Step-by-step re-implementation of Karpathy code for llama2.

1 XXX

1.1 Conventions

If T is a table of dimensions $n_1 \times n_2$, then T_i is a vector of length n_2

1.2 Glossary

Abbreviation	C code	Comments
n_D	dim	
n_{HD}	hidden_dim	
n_L	n_layers	
n_H	n_heads	
n_{HKV}	n_kv_heads	
n_V	vocab_size	
n_S	seq_len	
s_H	head_size	
n'_D	ispc_dim	

Table 1: Scalars: Mapping math notation to C code

Abbreviation	C code	Dimensions	Comments
W_t	token embedding table	2	$n_V \times n_D$
W_{att}	rms_att_weight	2	$n_L \times n_D$
W_{ffn}	rms_ffn_weight	2	$n_L \times n_D$
W_q	w_q	3	$n_L \times n_D \times (n_H \times s_H)$
W_k	w_k	3	$n_L \times n_D \times (n_HKV \times s_H)$
W_v	w_v	3	$n_L \times n_D \times (n_HKV \times s_H)$
W_o	w_o	3	$n_L \times (n_H \times s_H) \times n_D$

Table 2: Weights: Mapping math notation to C code

2 Utilities

```

 $x \leftarrow W_t[t]$ 

for each layer  $l$  do
    attention rmsnorm
     $x_b \leftarrow \text{rmsnorm}(x, W_{ra}[l])$ 
    locate key and value in cache
     $x_{kc} \leftarrow (S_{kc}[l])[p]$ 
     $x_{vc} \leftarrow (S_{vc}[l])[p]$ 
     $XXX \leftarrow \text{matmul}(x_{kc}, x_b, W_{wk}[l], )$  SOME JUNK
endfor

```

Figure 1: Pseudo-code for forward

3 Forward

Arguments in Table 3

Argument	Type	Comments
T	Transformer	pointer
t	integer	token $0 \leq t < n_V$
p	integer	pos

Table 3: Arguments for forward

3.1 rmsnorm

Input is in Table 4. Output is o , float vector of length n

Argument	Type
x	float vector of length n
w	float vector of length n

Table 4: Arguments for rmsnorm

$$1. \alpha = ((\sum_i x_i^2)/n) + \epsilon$$

$$2. o_i \leftarrow \frac{w_i \times x_i}{\alpha}$$

3.2 softmax

Arguments in Table 5

Argument	Type
x	float vector of length n
n	integer

Table 5: Arguments for softmax

1. $m = \max_{i=0}^{i=n-1} x_i$
2. $\forall_{i=0}^{i=n-1} x_i \leftarrow e^{x_i - m}$
3. $s = \sum_{i=0}^{i=n-1} x_i$
4. $\forall_{i=0}^{i=n-1} x_i \leftarrow \frac{x_i}{s}$