**Algorithm:** $P \leftarrow$ `CodeLlama`$(x|\theta)$

---

**Input:** $x \in V^*$, a sequence of token IDs.

**Output:** $P \in [0,1]^{N_V \times \text{length}(x)}$, where $P_t$ represents the conditional distribution $p(x_{t+1}| \, x_{1:t})$.

**Hyperparameters:** $N_V$, $\theta_{\max}$, $L$, $H$, $d_e$, $d_{\text{mlp}}$, MLP hidden dimension.

**Parameters:** $\theta$ includes all the following parameters:

$\quad W_e \in \mathbb{R}^{d_e \times N_V}$ $\qquad$ token embeddings matrix.

$\quad W_p \in \mathbb{R}^{d_e \times \theta_{\max}}$ $\qquad$ positional embeddings matrix.

$\quad W_l, W_{e2d_l}$ $\qquad\qquad$ attention and MLP weights.

$\quad$ For $l \in [L]$ $\qquad\qquad$ See paper for details.

$\quad \gamma, \beta \in \mathbb{R}^{d_e}$ $\qquad\qquad$ output normalization parameters.

$\quad W_u \in \mathbb{R}^{N_V \times d_e}$ $\qquad\quad$ unembedding matrix.

1 $\quad \theta \leftarrow \text{length}(x)$

2 $\quad$ for $t \in \theta : e_t \leftarrow W_e[:, x[t]] \; + \; W_p[:, t]$

3 $\quad X \leftarrow [e_1, e_2, \ldots e_\theta]$

4 $\quad$ **for** $l = 1, 2, \ldots, L$ **do**

5 $\qquad X \leftarrow$ `layer_norm`$(X[:, t] \, | \, \gamma_l^1, \beta_l^1)$

6 $\qquad X \leftarrow X + $ `MHAttention`$(\tilde{X} \, | \, W_l, \text{Mask}[t, t'] = [[t \leq t']])$

7 $\qquad X \leftarrow$ `layer_norm`$(X \, | \, \gamma, \beta)$

8 $\qquad X \leftarrow X + \text{MLP}(X \, | \, W_l)$

9 $\quad X \leftarrow$ `layer_norm`$(X \, | \, \gamma, \beta)$

10 $\quad$ **end**

11 $\quad$ for $t \in [\theta] : P_t \leftarrow \text{softmax}(W_u X[:, t])$

12 $\quad$ **return** $P$

This implements an autoregressive decoder-only transformer similar to GPT with causal self-attention masking. It processes the input sequence, applies multiple transformer layers with causal self-attention and MLP blocks, and produces a conditional distribution over the next token at each position. The specifics like layer normalization and multi-head attention are as defined in the "Formal Algorithms for Transformers" article.