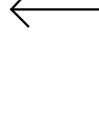
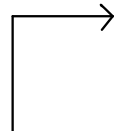


\times
 N Layers

Vanilla Encoder Attention Transformer

Vanilla Encoder Attention Transformer

Fuse And Pool CLS Token
 +
 Feed Forward Classification Head



$X_T \in \mathbb{R}^{T_l \times d}$



$X_C \in \mathbb{R}^{C_l \times d}$

$X_T \in \mathbb{R}^{T_l \times d}$

+

E_r

+

E_c

$cell_1$

$cell_2$

$cell_{T_l}$

E_{r_1}

E_{r_1}

$E_{r_{T_r}}$

E_{c_1}

E_{c_2}

$E_{c_{T_c}}$

$X_C \in \mathbb{R}^{C_l \times d}$

+

P_C

C_1

C_2

C_{C_l}

P_1

P_2

P_{C_l}