**Output Probability**

**Feed Forward**

**Norm**

**Cross-Attention & Gating Network**

**Multi-Head Self-Attention**

**Norm**

**Input Tokens**

**Transfomer Block**

**Gate**

Softmax

Softmax

$W_q$

$W_{v(i)}$

$W_{k(j)}$

**Norm**

**Norm**

**Norm**

Patent title & abstract

Legal status

BT + Time

Labels

**Cross-Attention & Gating Network**