

Output:

$p(s_1)$

$p(s_2)$

...

$p(s_n)$

Softmax

Generator

Layer normalization

Feed Forward network (w2)

Feed Forward network (w1)  
with activation

Layer normalization

out

self-attention

key

value

query

Multi-head attention

positional encoding

Input:

trace 1

[start]

[m]

$s_2$

...

$s_{k-2}$

$s_{k-1}$

[stop]

trace 2

[start]

$s_1$

$s_2$

...

[m]

$s_{k-1}$

[stop]