# Decoder Layer

## Multihead Attention Layer

vector    W    output     dropout 0.1     Norm Layer

...  ...

d_model

$$x'_i = \frac{x_i - \overline{X}}{S}$$

## Multihead Attention Layer

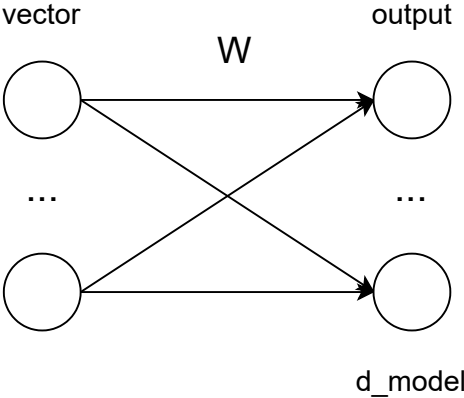vector    W    output     dropout 0.1     Norm Layer

...  ...

d_model

$$x'_i = \frac{x_i - \overline{X}}{S}$$

## Feed Forward Network

input    W    output     dropout 0.1     Norm Layer

...  ...  ...

d_model      d_model

$$x'_i = \frac{x_i - \overline{X}}{S}$$

output