

Detailed Analysis of the Encoder

◆ Single Encoder Layer

1

Multi-Head Self-Attention

Captures relationships between all positions

+ Residual → Layer Norm

2

Feed-Forward Network

Position-wise transformation

+ Residual → Layer Norm



Residual Connections

Around each sub-layer to help gradient flow



Layer Normalization

Applied after each sub-layer for stability



Same Structure

All layers share identical structure but have different parameters

Output dimension preserved throughout:

$d_{model} = 512$



This structure is repeated 6 times (stacked layers)

Each layer has different parameters