

Input sequence

Token embedding

+

Positional encoding

Encoder block

Multi-head attention

Add + Norm

Feed-forward network

Add + Norm

$\times N$
(stacked)

Contextual representations

