

Output Representations

Residual Stream

Input Features

N

Feedforward Network

Layer Normalization

Multi-head Self Attention

Layer Normalization

