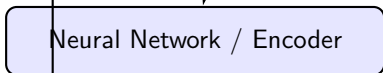


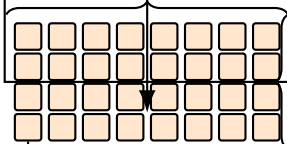
Hidden States $\mathbf{H} [B, L, D]$



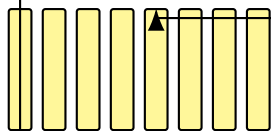
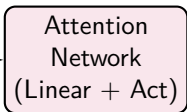
[batch, seq_len]



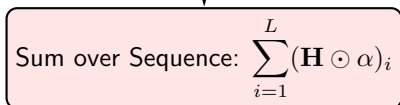
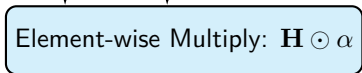
L positions



D
dim



Attention Weights $\alpha [B, L, 1]$



Pooled Output $\mathbf{h}_{pool} [B, D]$