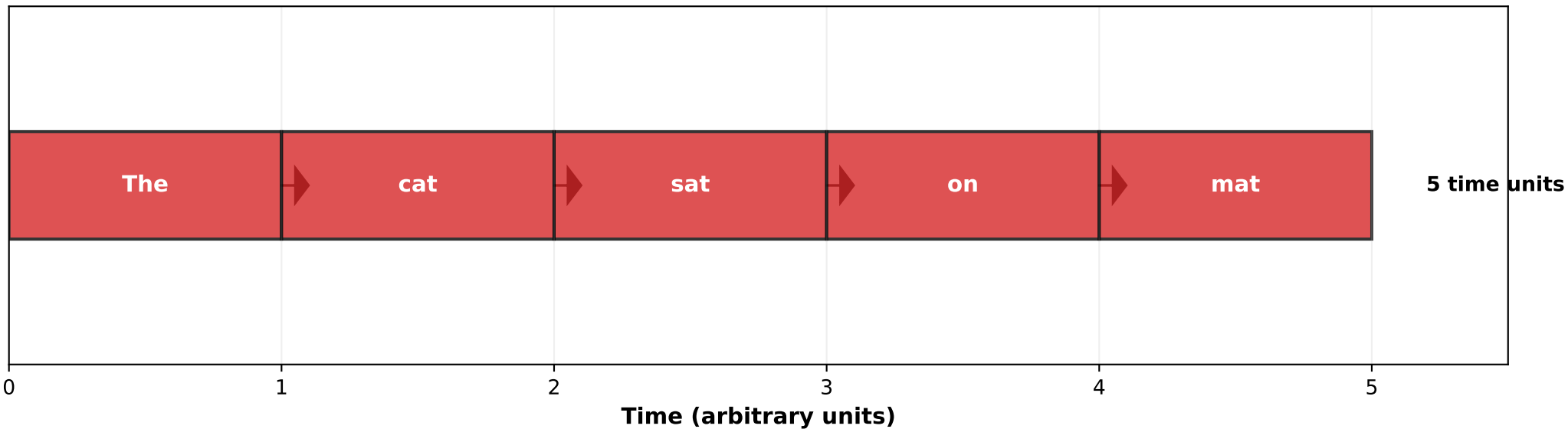


Sequential vs Parallel: The Bottleneck

RNN: Sequential Processing (One Word at a Time)



Transformer: Parallel Processing (All Words Simultaneously)

