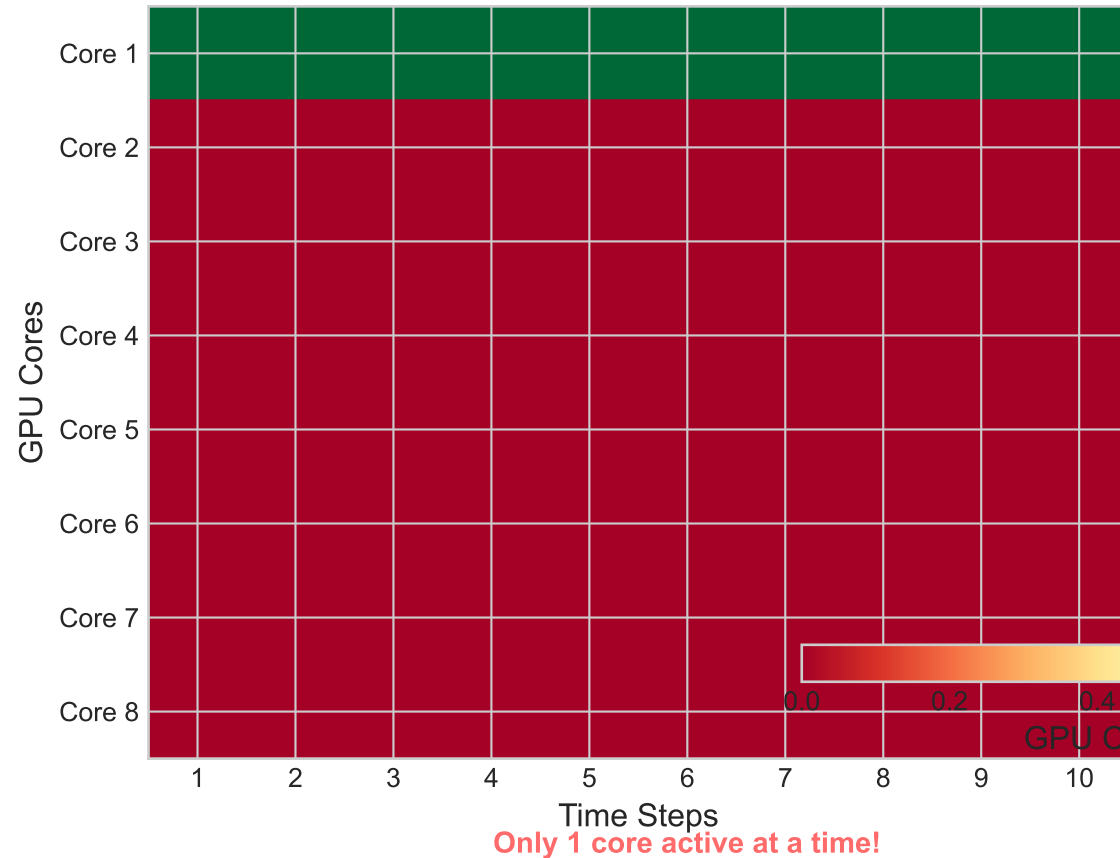


# GPU Efficiency: Why Transformers Train Faster

**RNN: Sequential Processing**  
(Low GPU Utilization)



**Transformer: Parallel Processing**  
(Full GPU Utilization)

