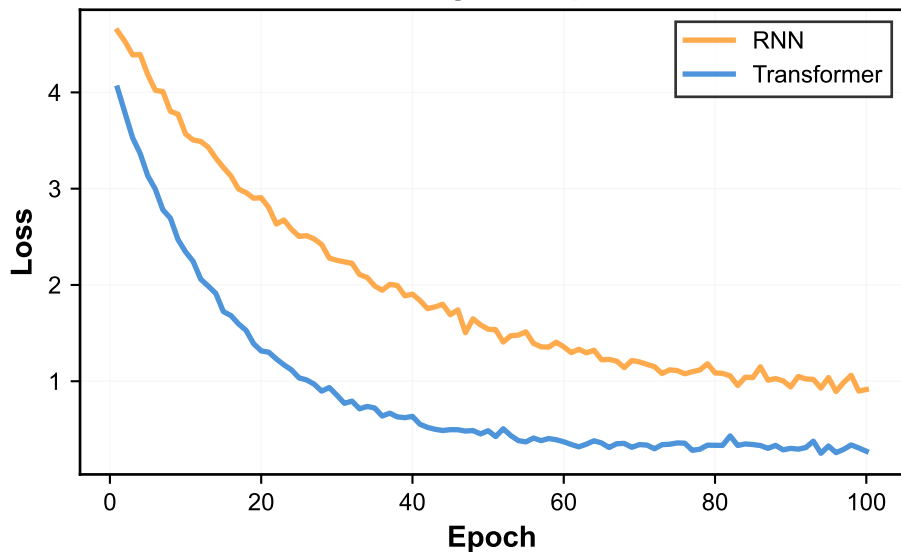
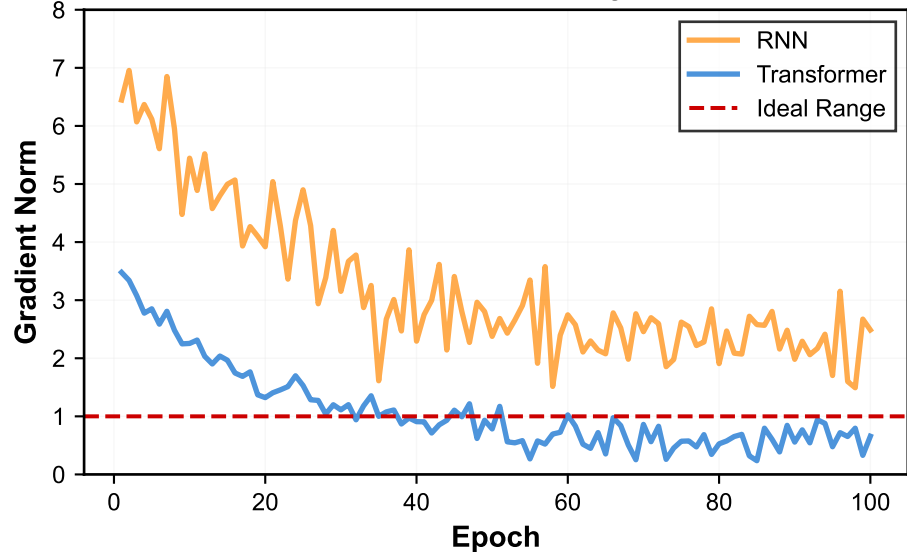


Training Dynamics: RNN vs Transformer

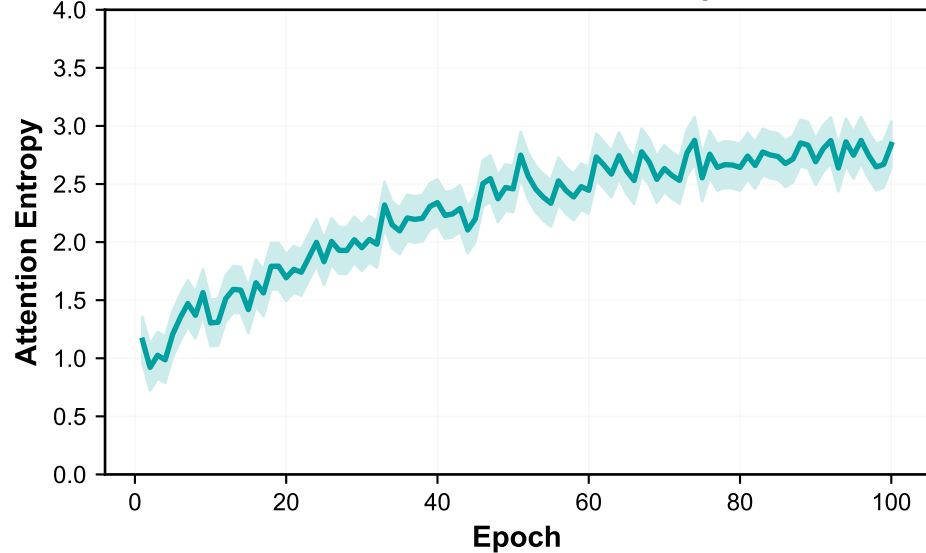
Convergence Speed



Gradient Stability



Attention Pattern Diversity



Final Performance

