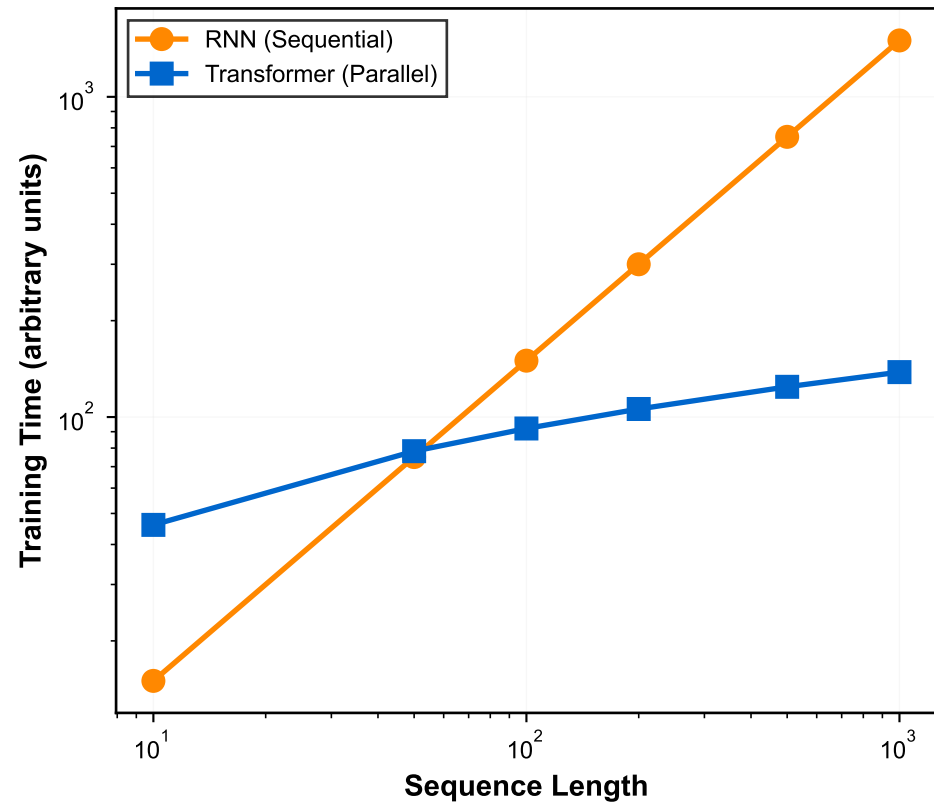


Why Transformers Train Faster

Training Time vs Sequence Length



GPU Efficiency: Parallel vs Sequential

