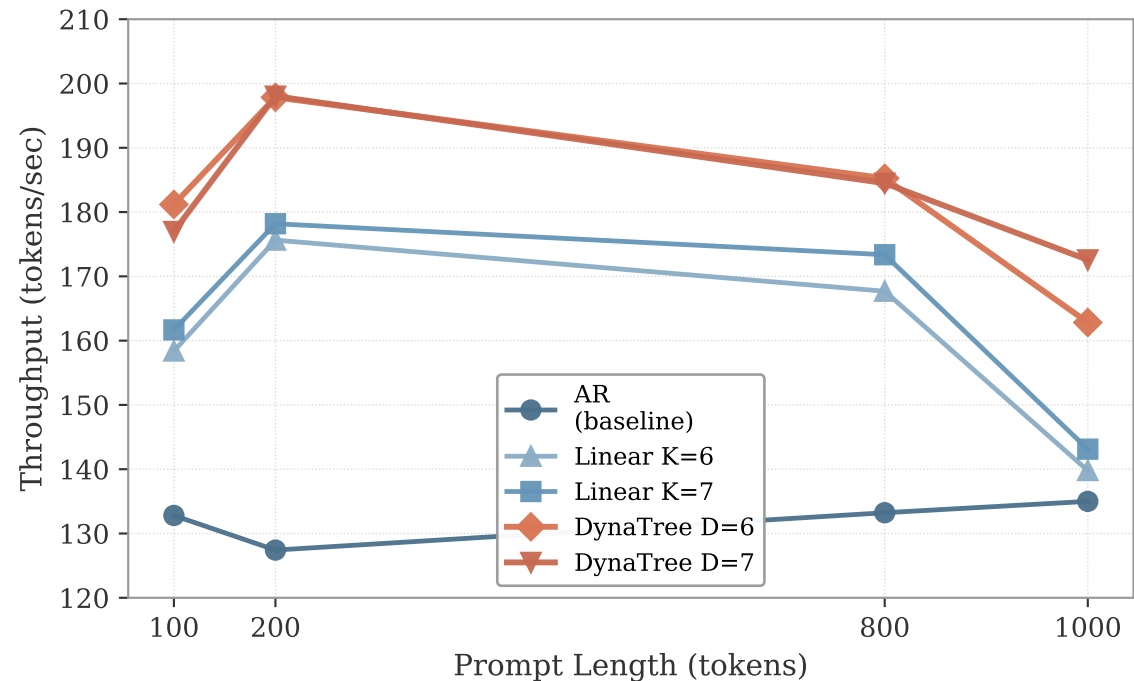


(a) Throughput vs. Prompt Length



(b) Speedup vs. Prompt Length

