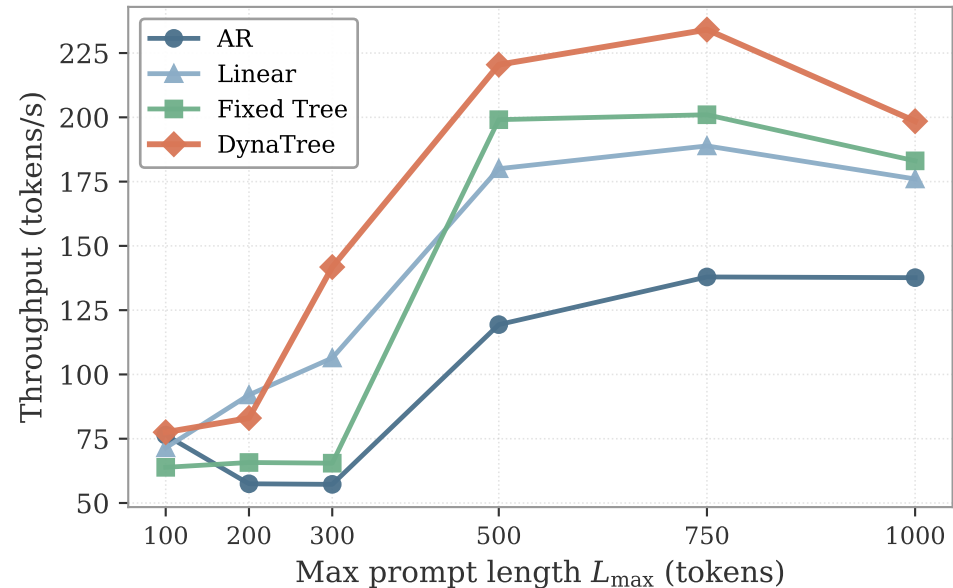(a) Throughput vs. prompt length

(b) Speedup vs. prompt length