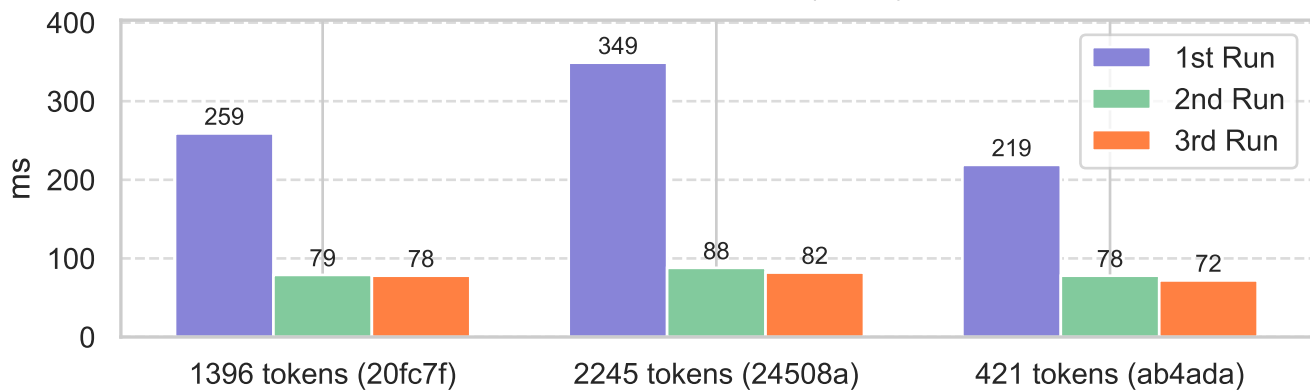
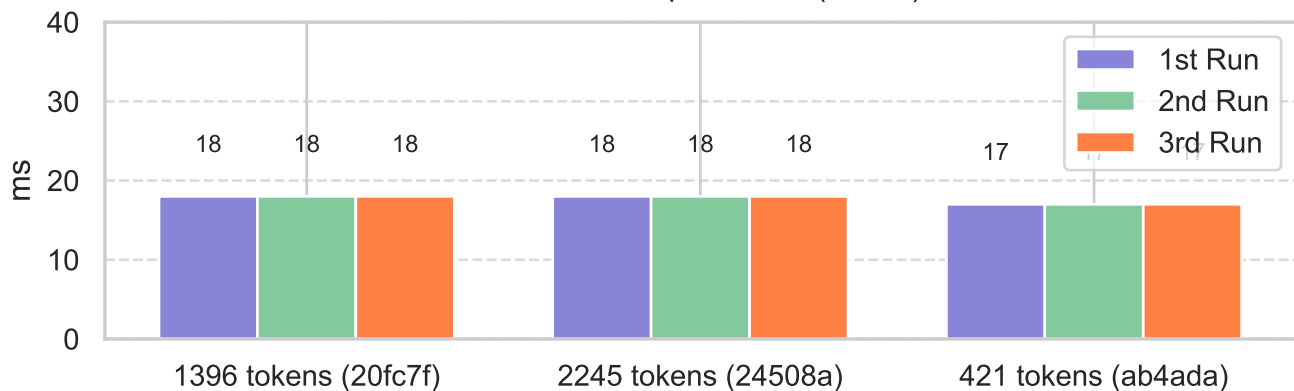


Time to First Token (TTFT)



Time Per Output Token (TPOT)



End-to-End Latency (E2E)

