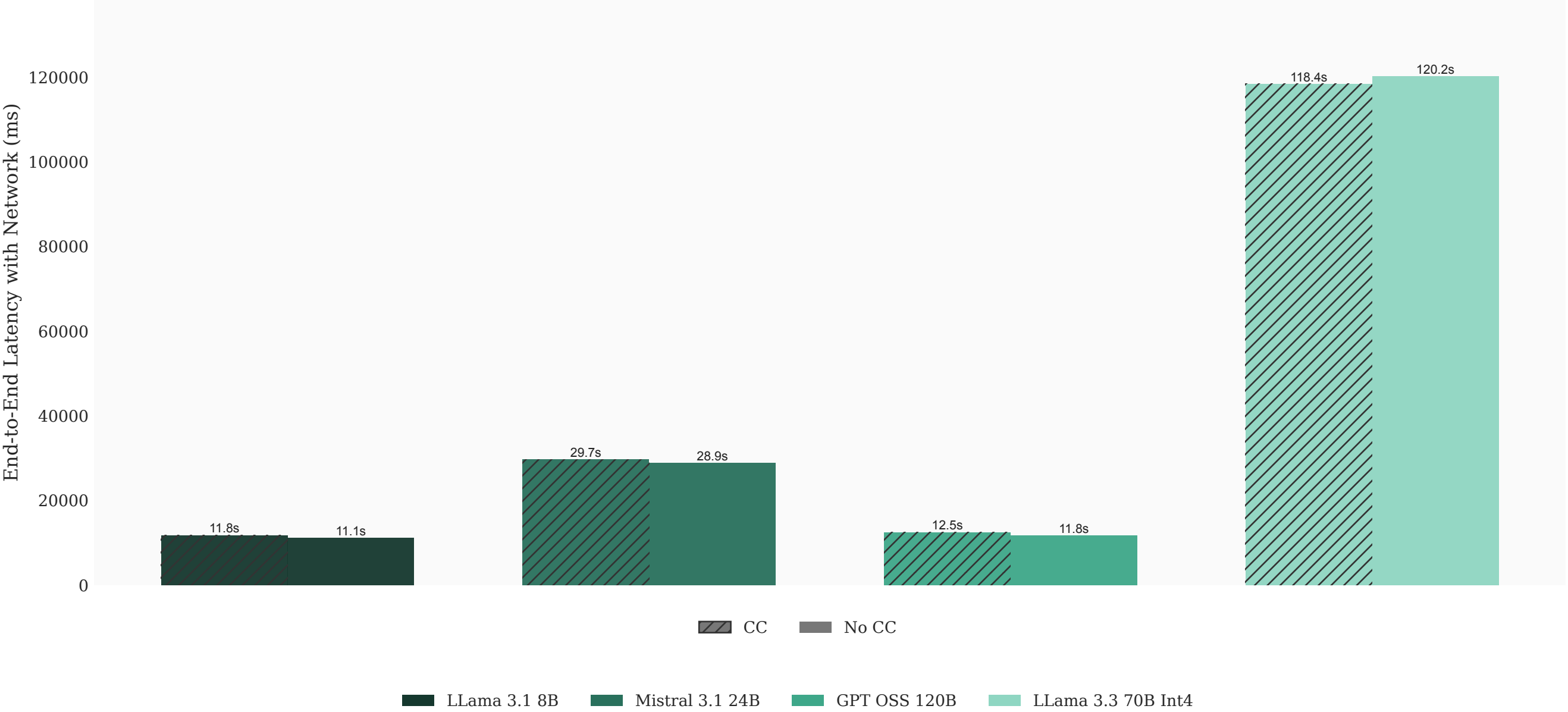


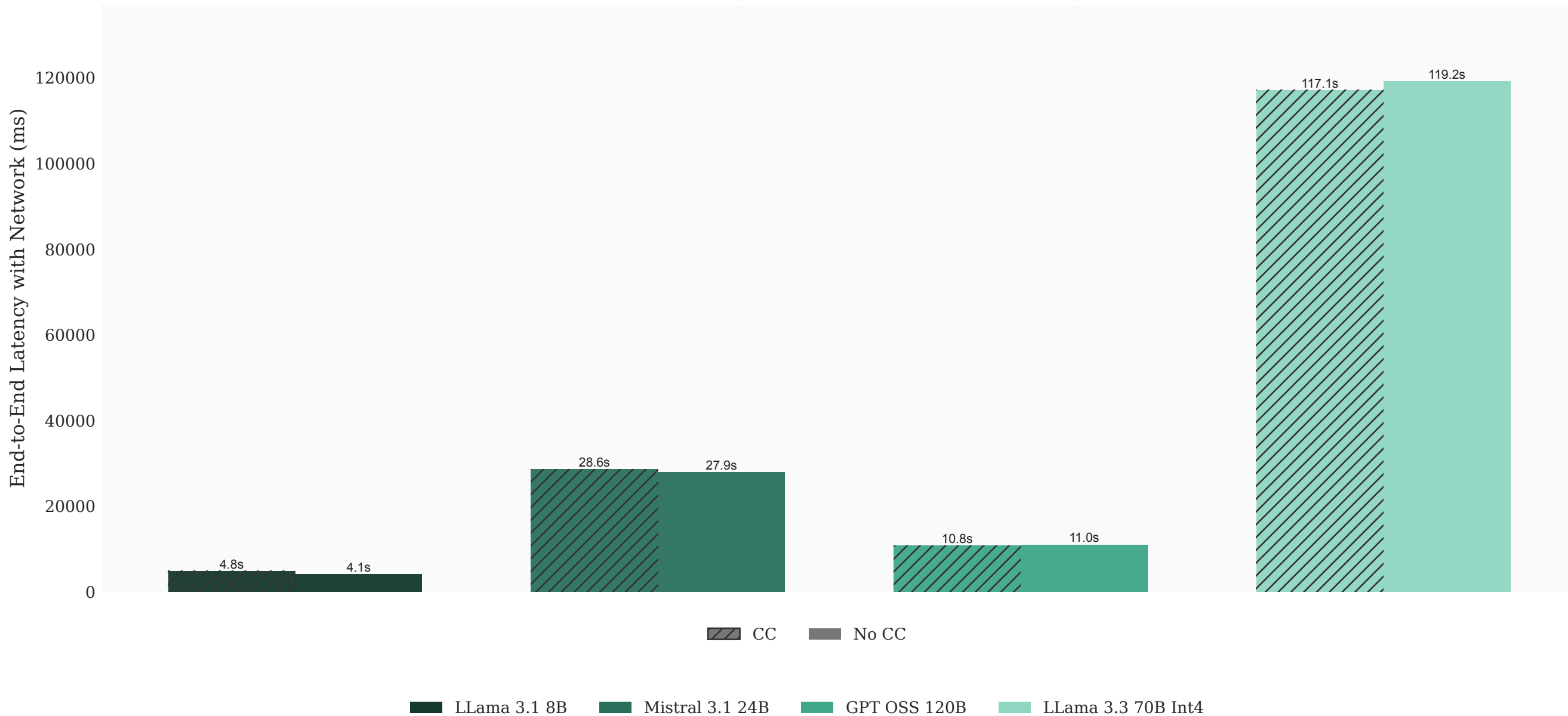
Random (1500  $\Rightarrow$  250) (Rate 100)

E2E Latency + 100ms Network Latency



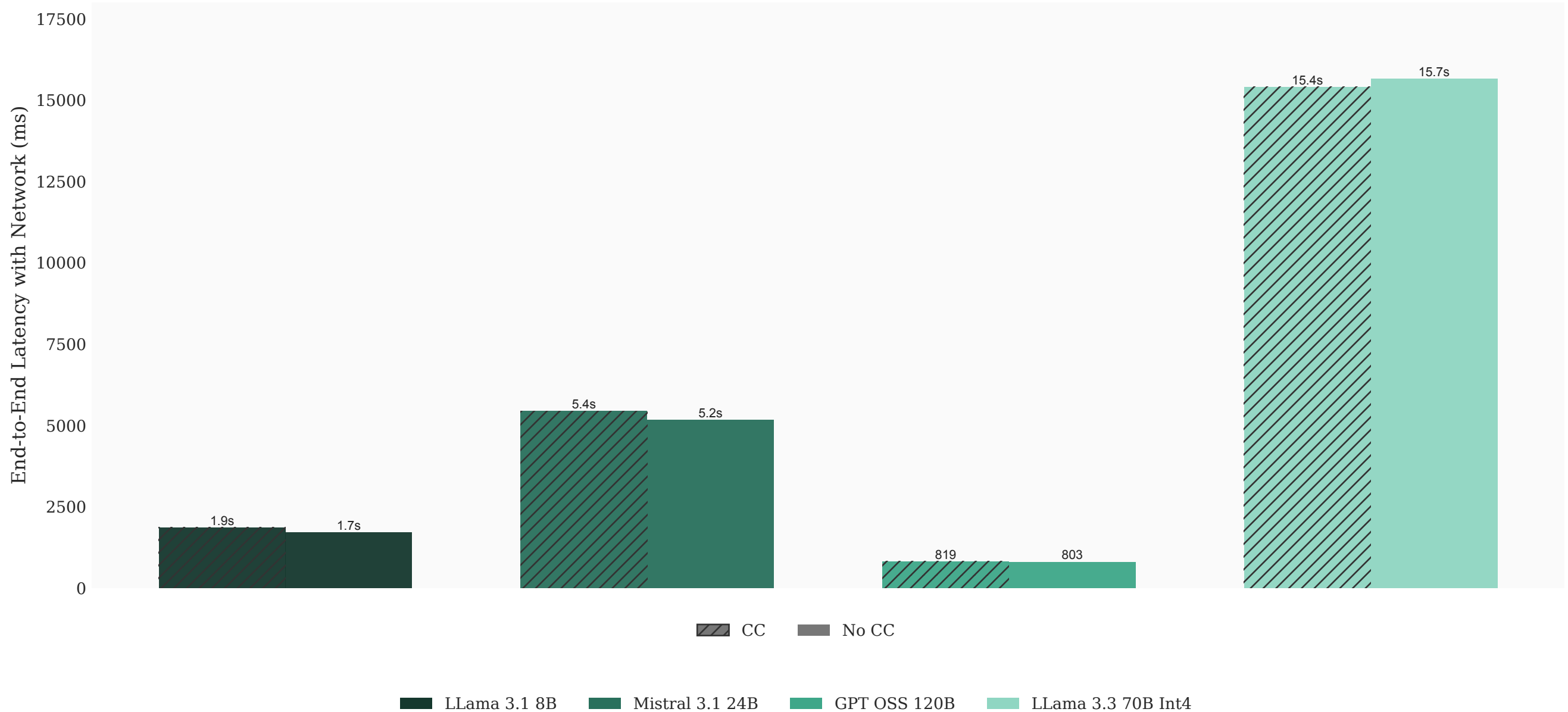
## Random (1500 $\Rightarrow$ 250) (Rate 50)

E2E Latency + 100ms Network Latency



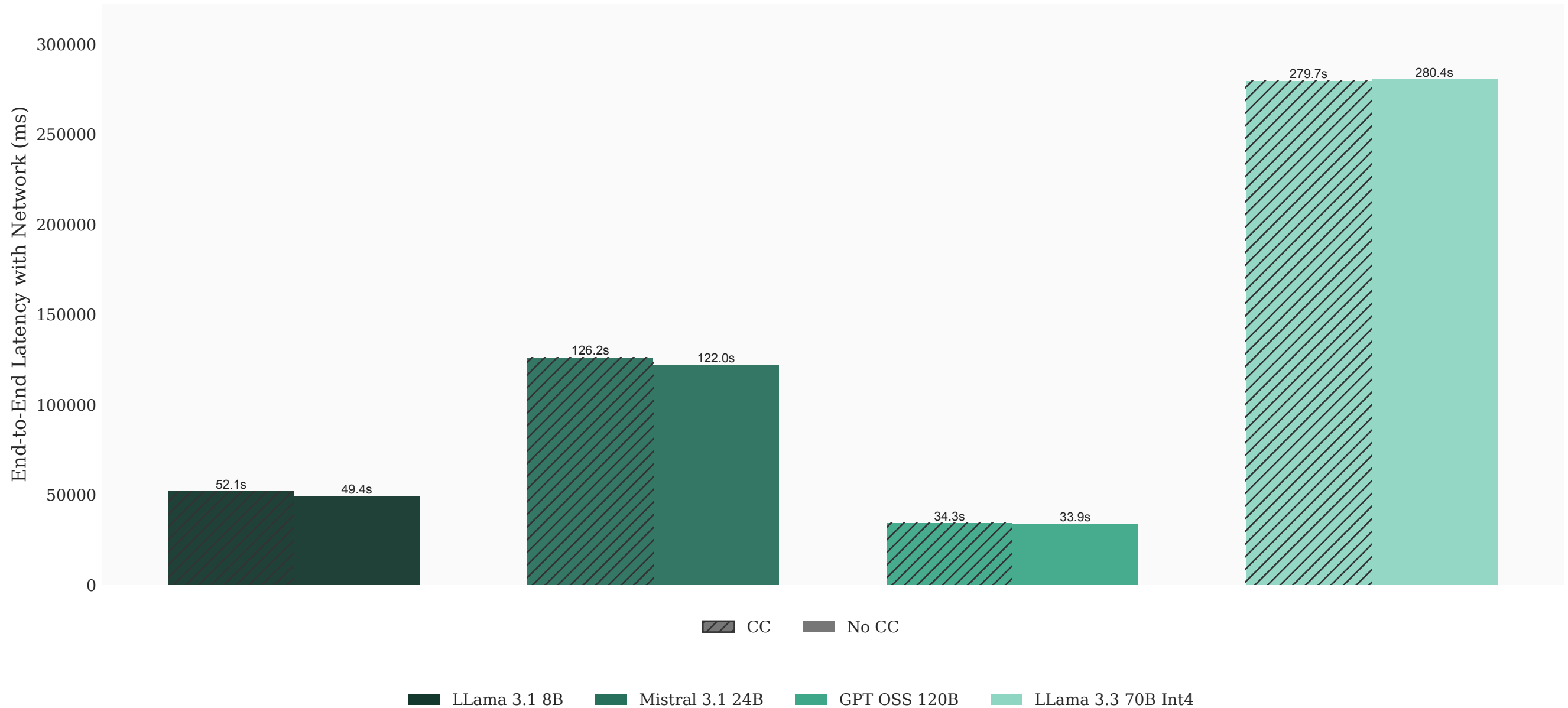
## Random (1500 $\Rightarrow$ 250) (Rate 1)

E2E Latency + 100ms Network Latency



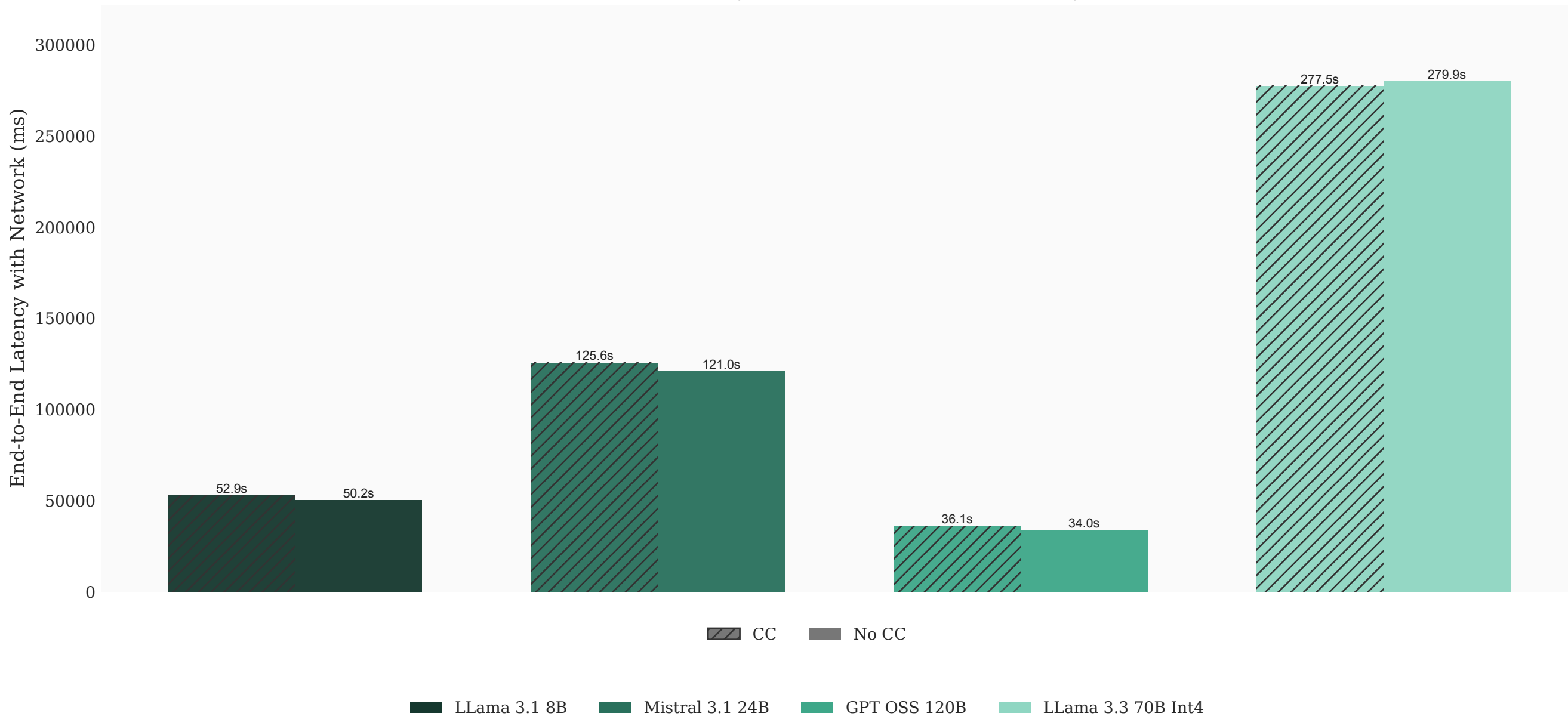
## Random (4000 $\Rightarrow$ 1000) (Rate 100)

E2E Latency + 100ms Network Latency



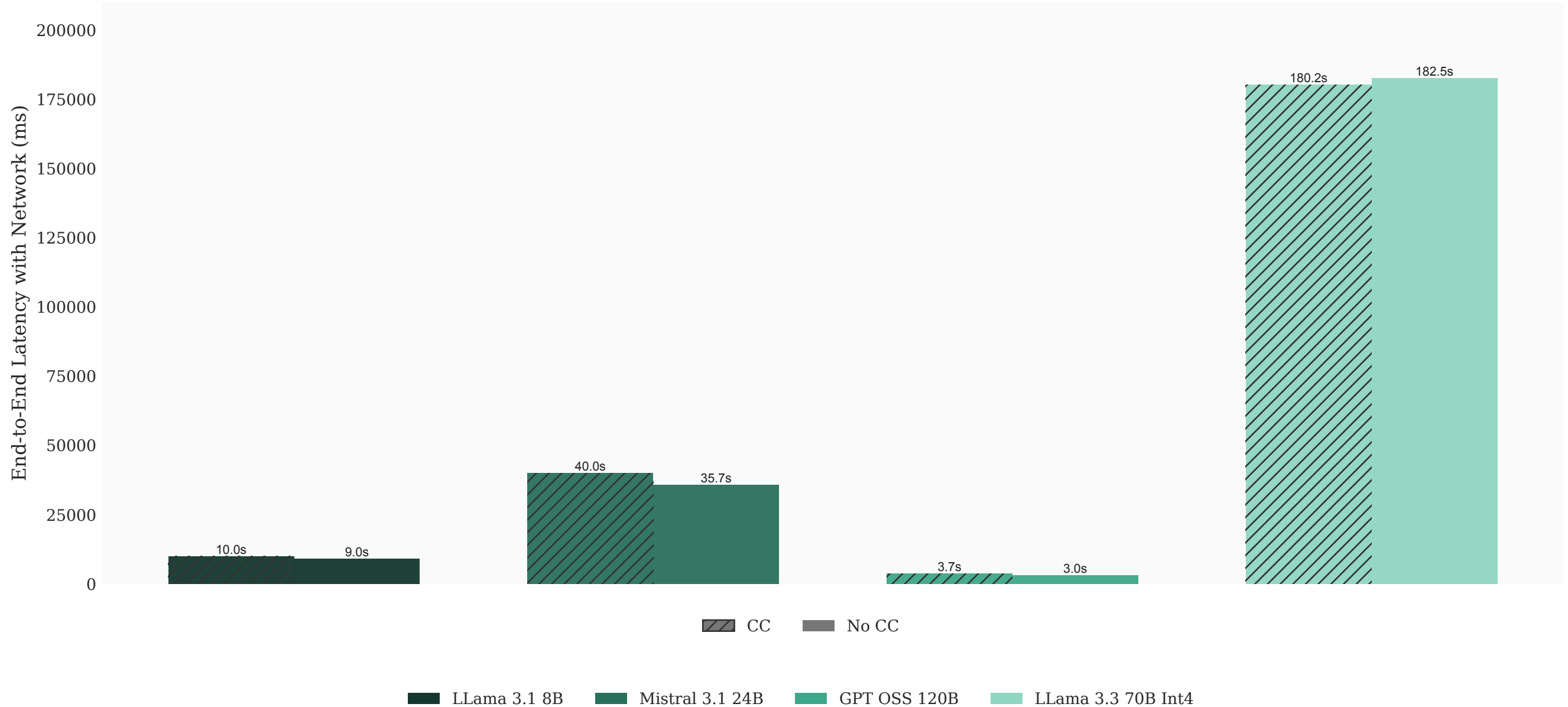
## Random (4000 $\Rightarrow$ 1000) (Rate 50)

E2E Latency + 100ms Network Latency



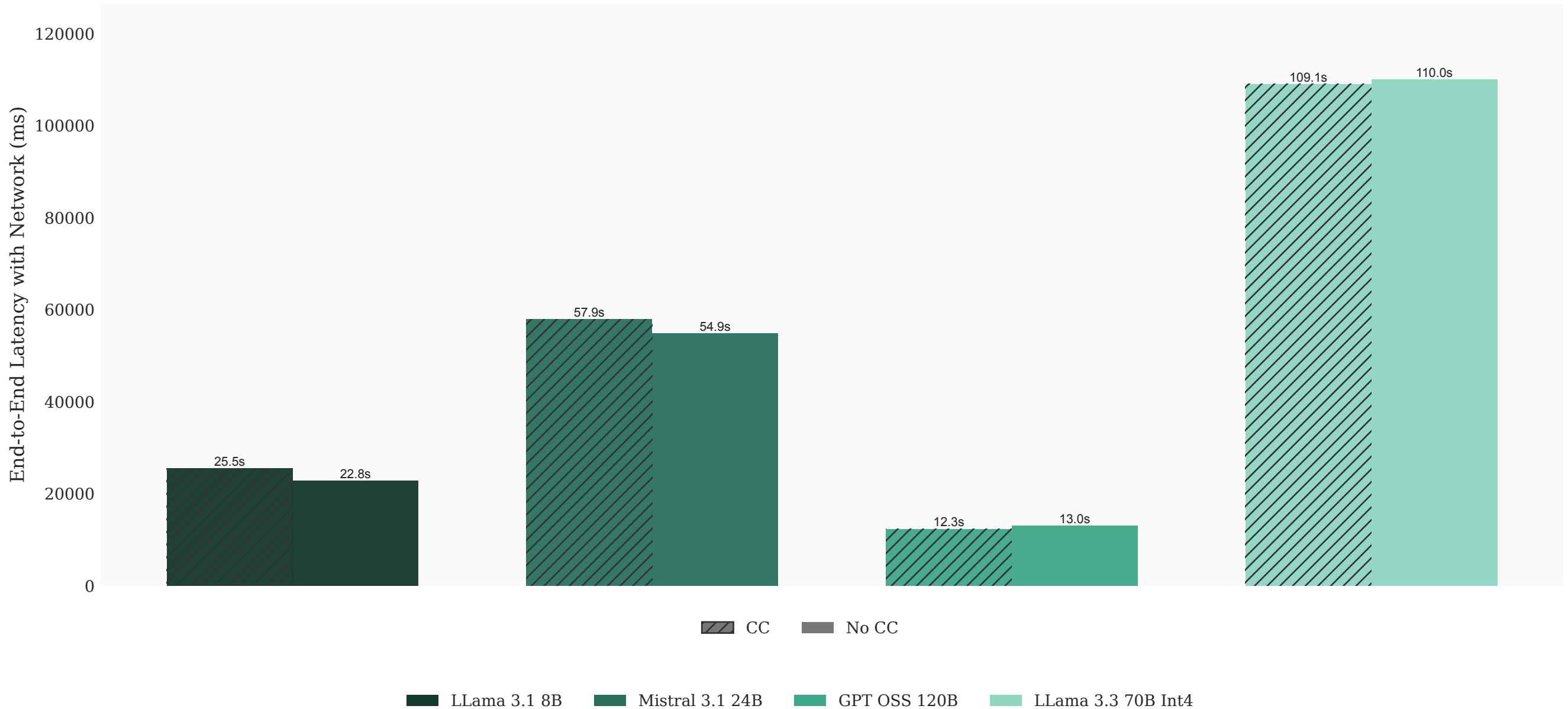
## Random (4000 $\Rightarrow$ 1000) (Rate 1)

E2E Latency + 100ms Network Latency



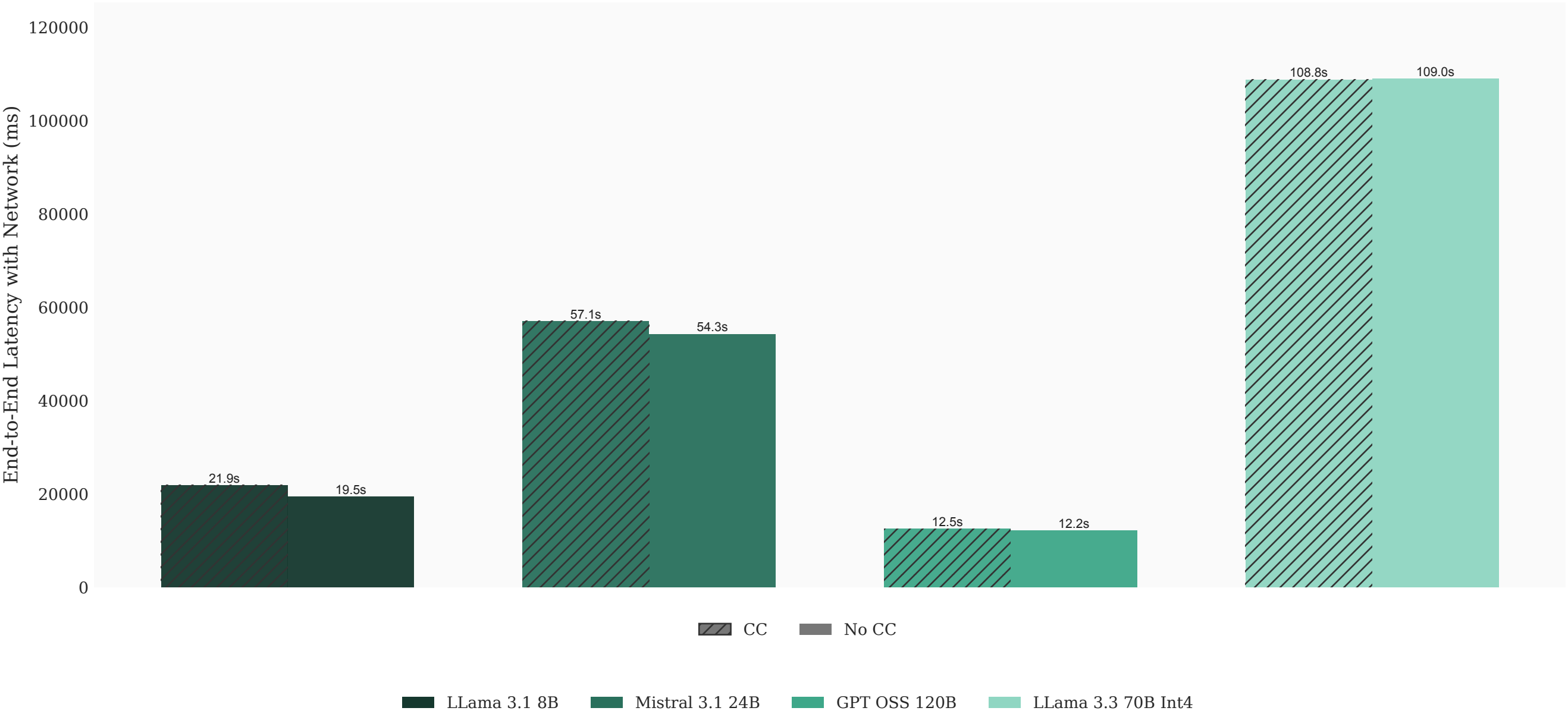
# Random (1000 $\Rightarrow$ 1000) (Rate 100)

E2E Latency + 100ms Network Latency



Random (1000  $\Rightarrow$  1000) (Rate 50)

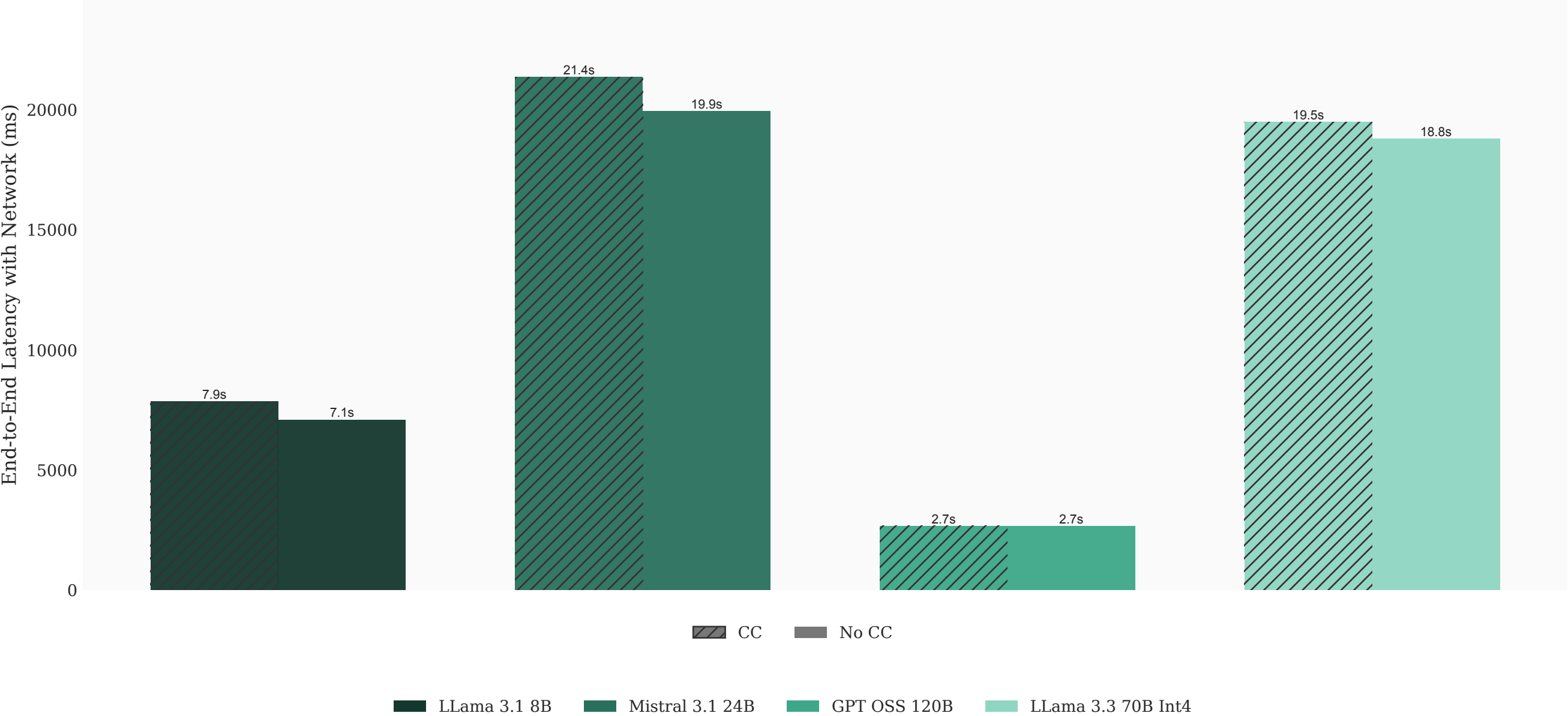
E2E Latency + 100ms Network Latency





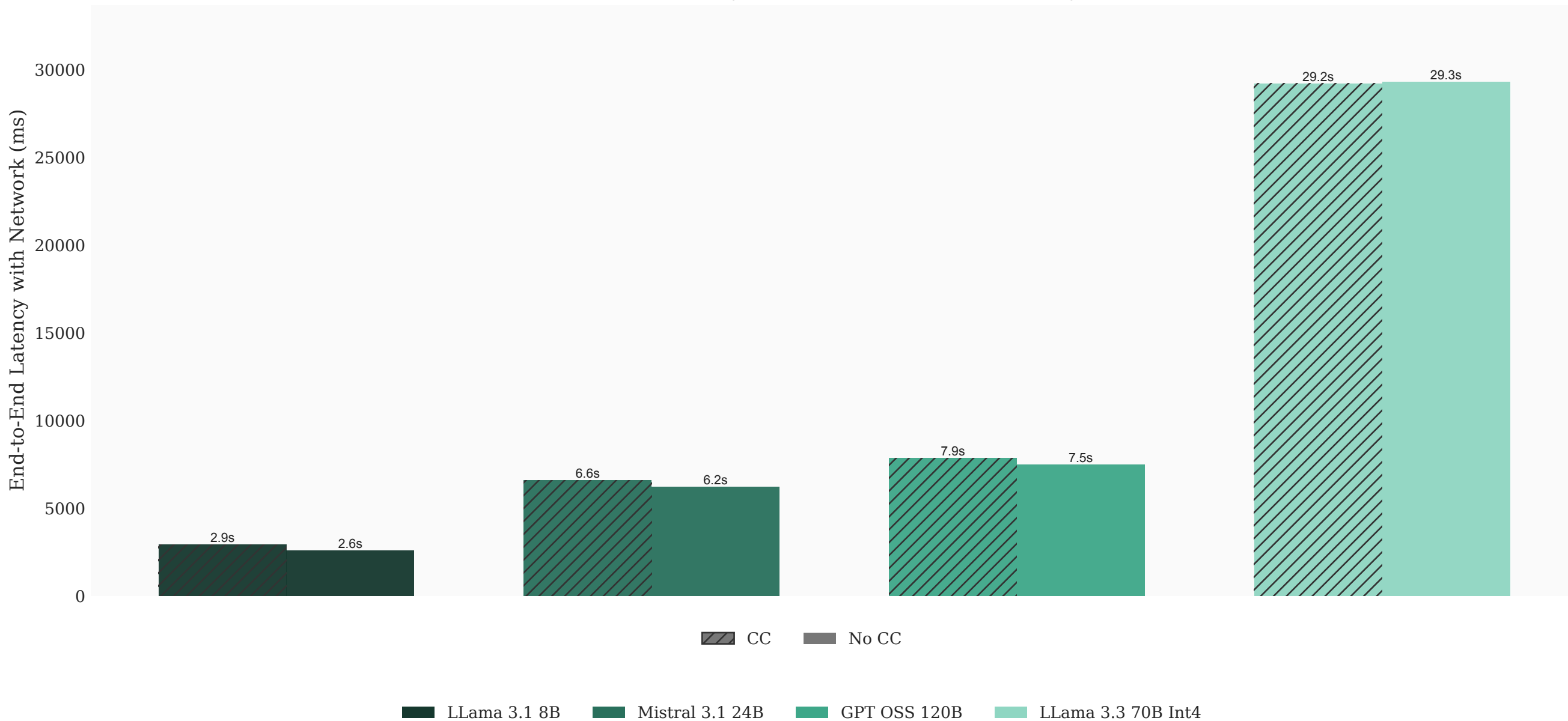
Random (1000 ⇒ 1000) (Rate 1)

E2E Latency + 100ms Network Latency



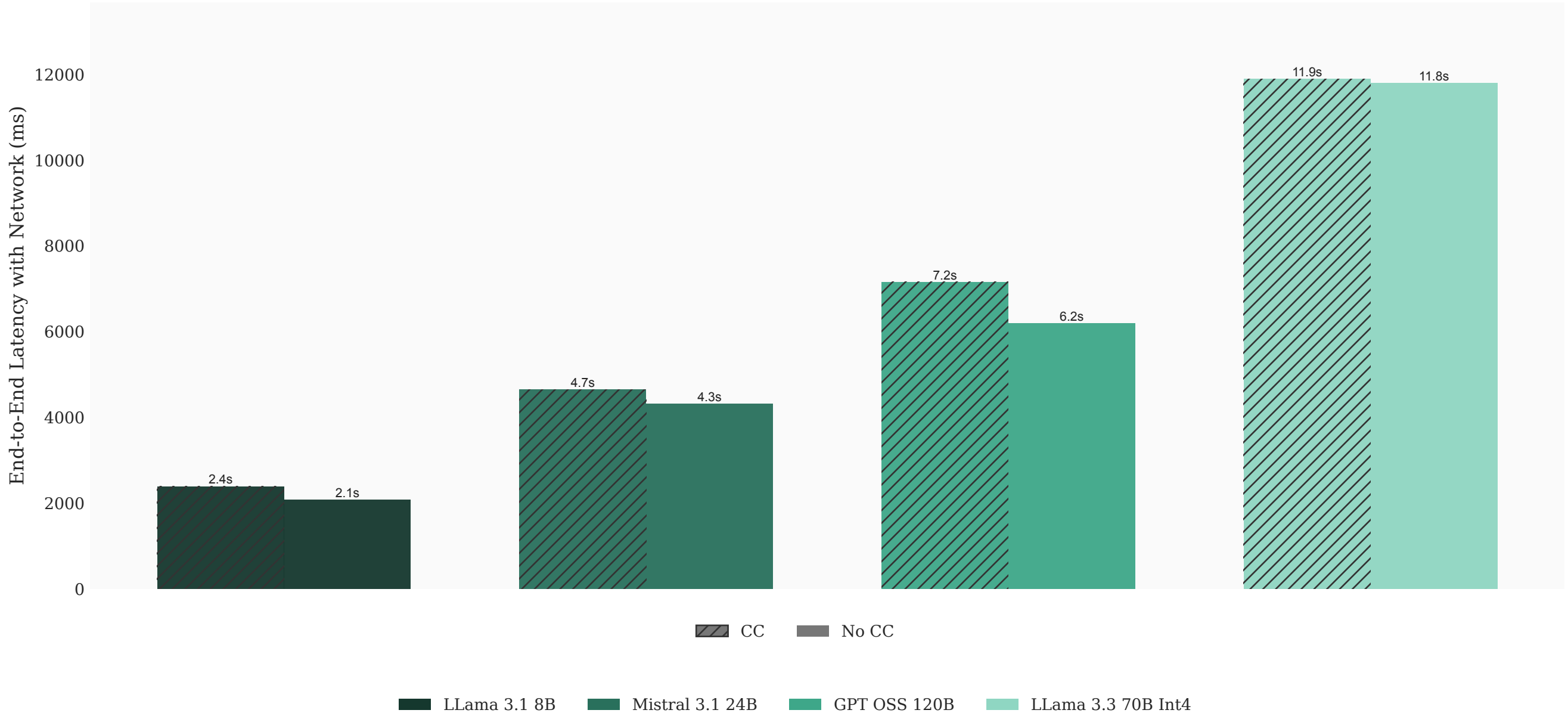
# ShareGPT (Rate 100)

E2E Latency + 100ms Network Latency



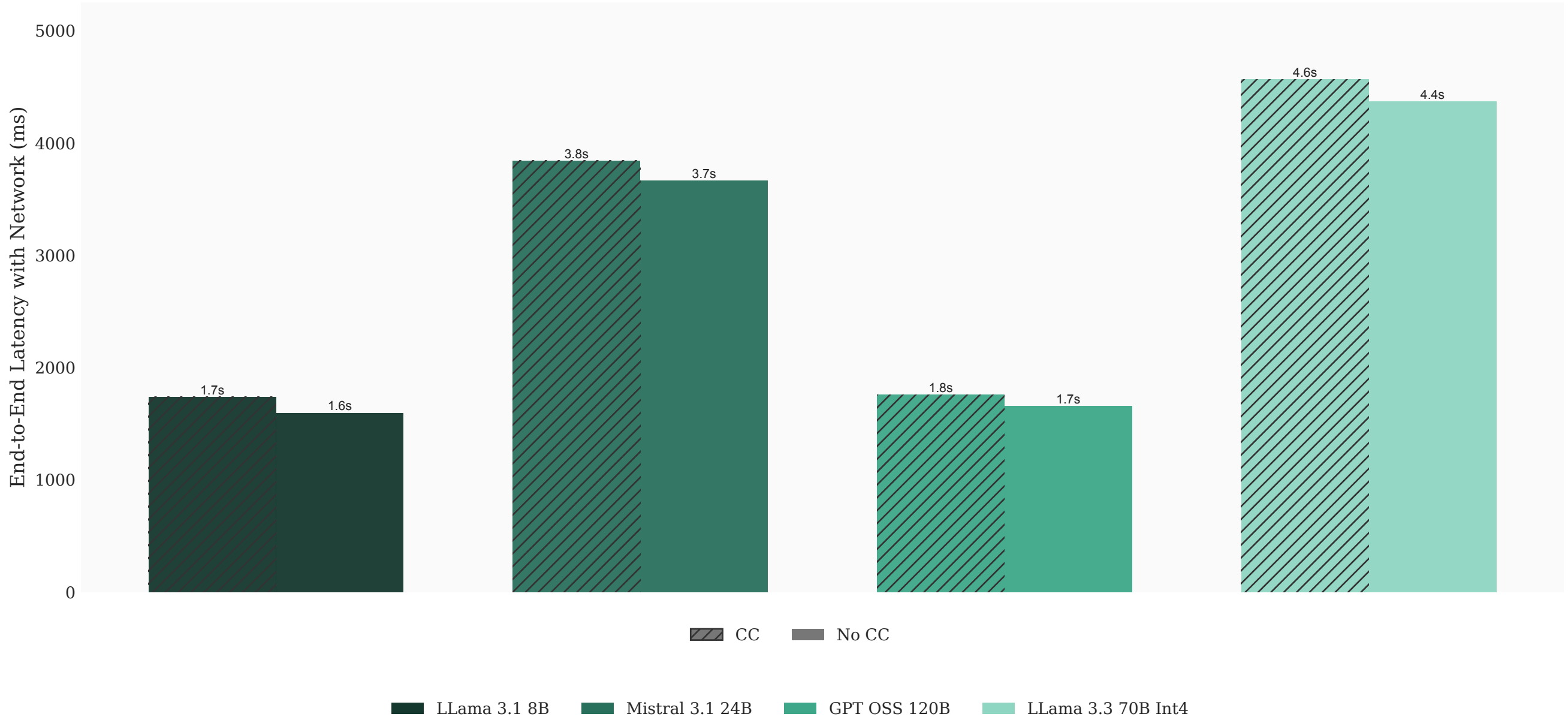
# ShareGPT (Rate 50)

E2E Latency + 100ms Network Latency



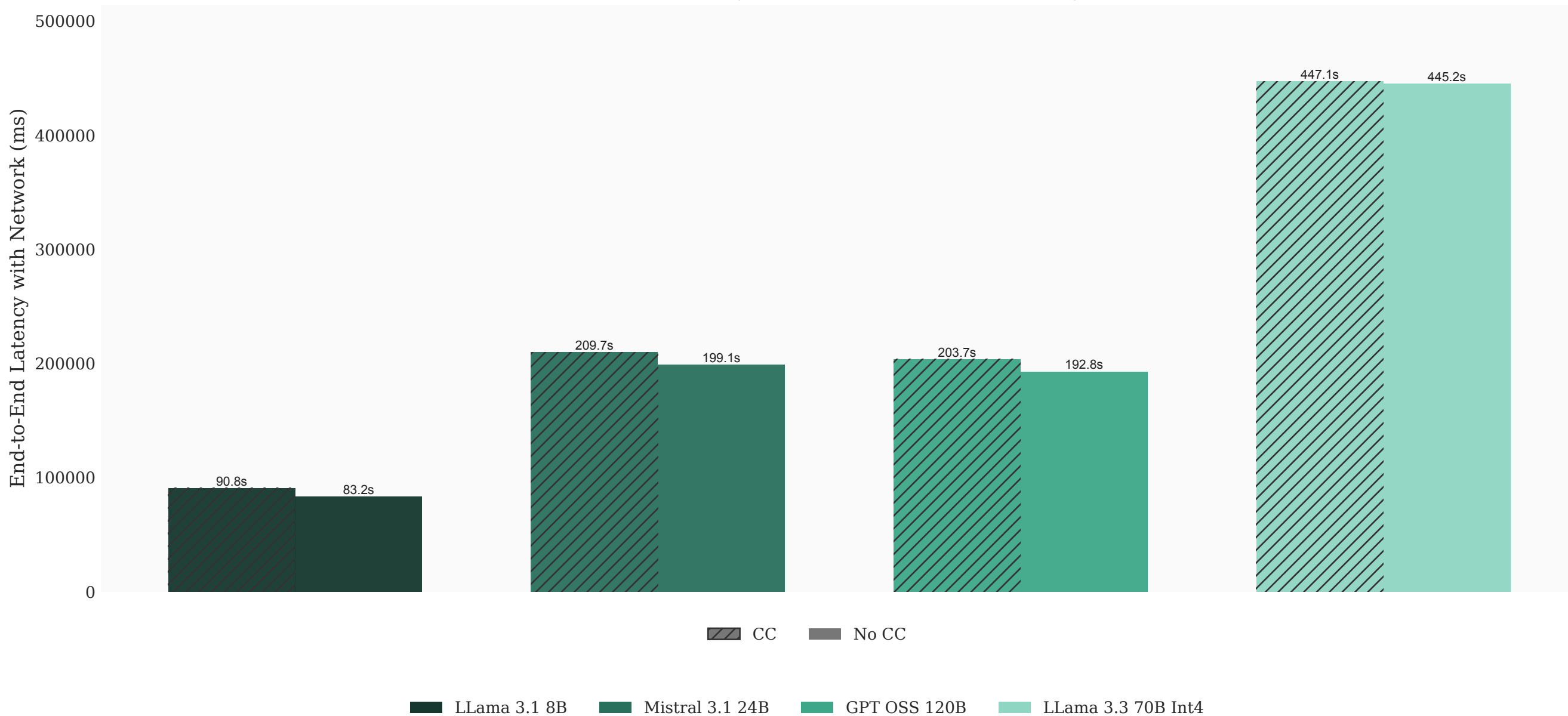
# ShareGPT (Rate 1)

E2E Latency + 100ms Network Latency



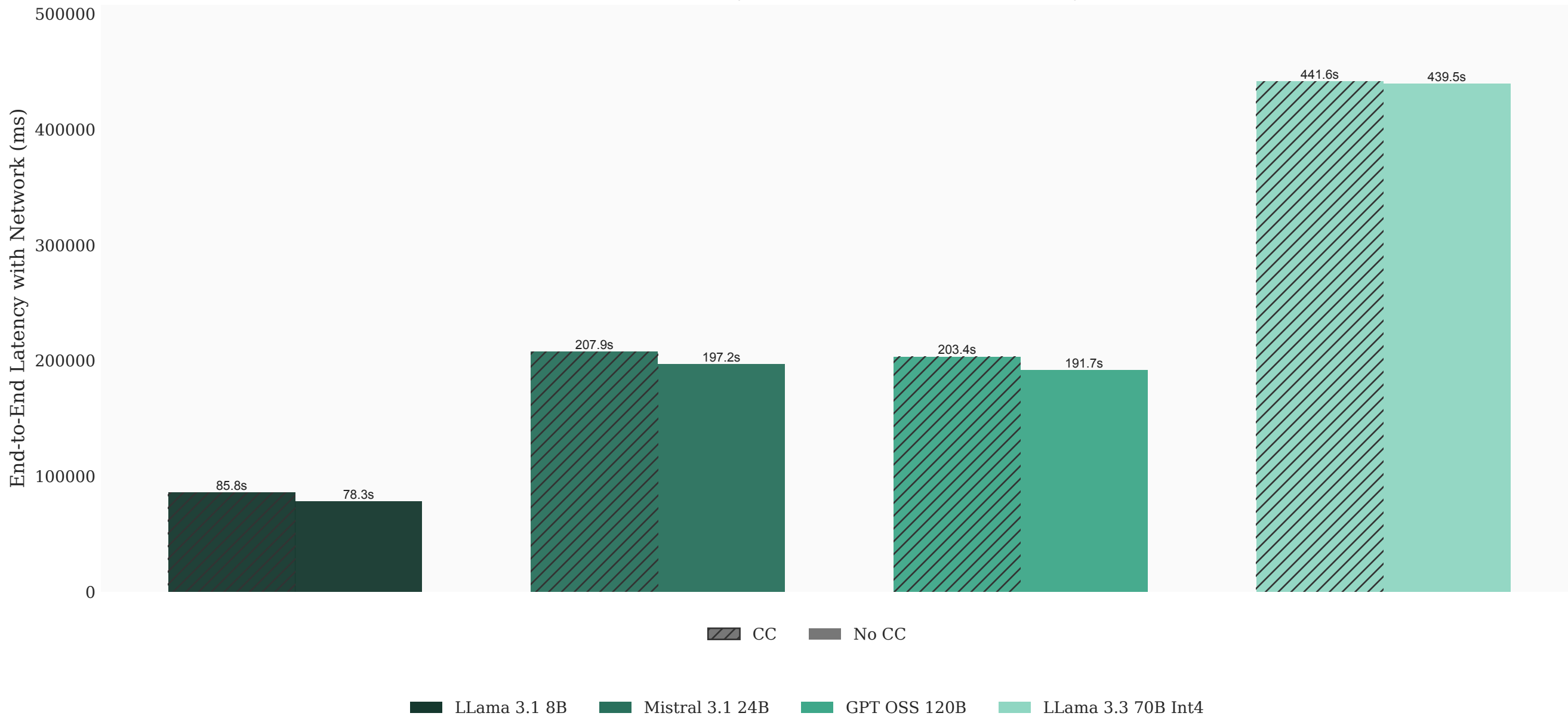
## Edit 10K Characters (Rate 100)

E2E Latency + 100ms Network Latency



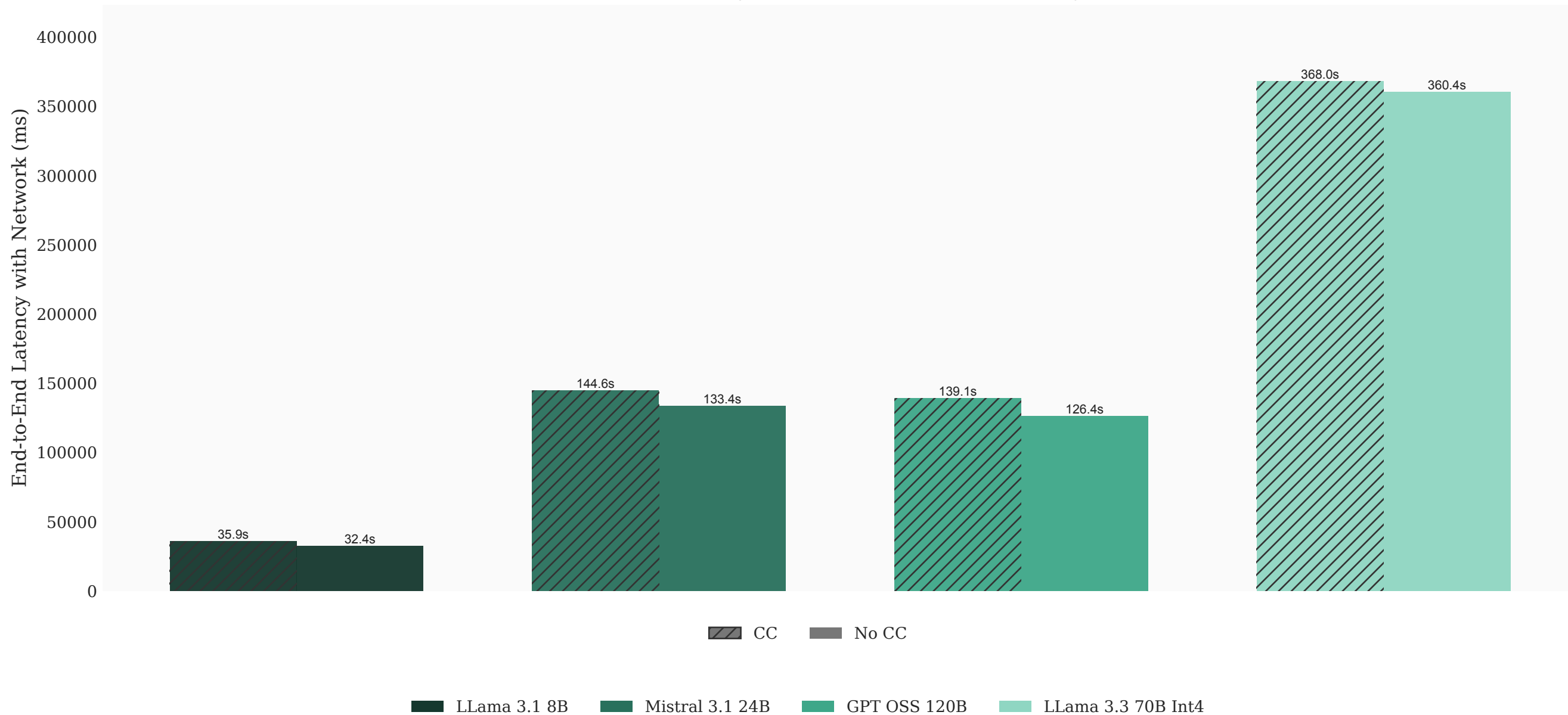
## Edit 10K Characters (Rate 50)

E2E Latency + 100ms Network Latency



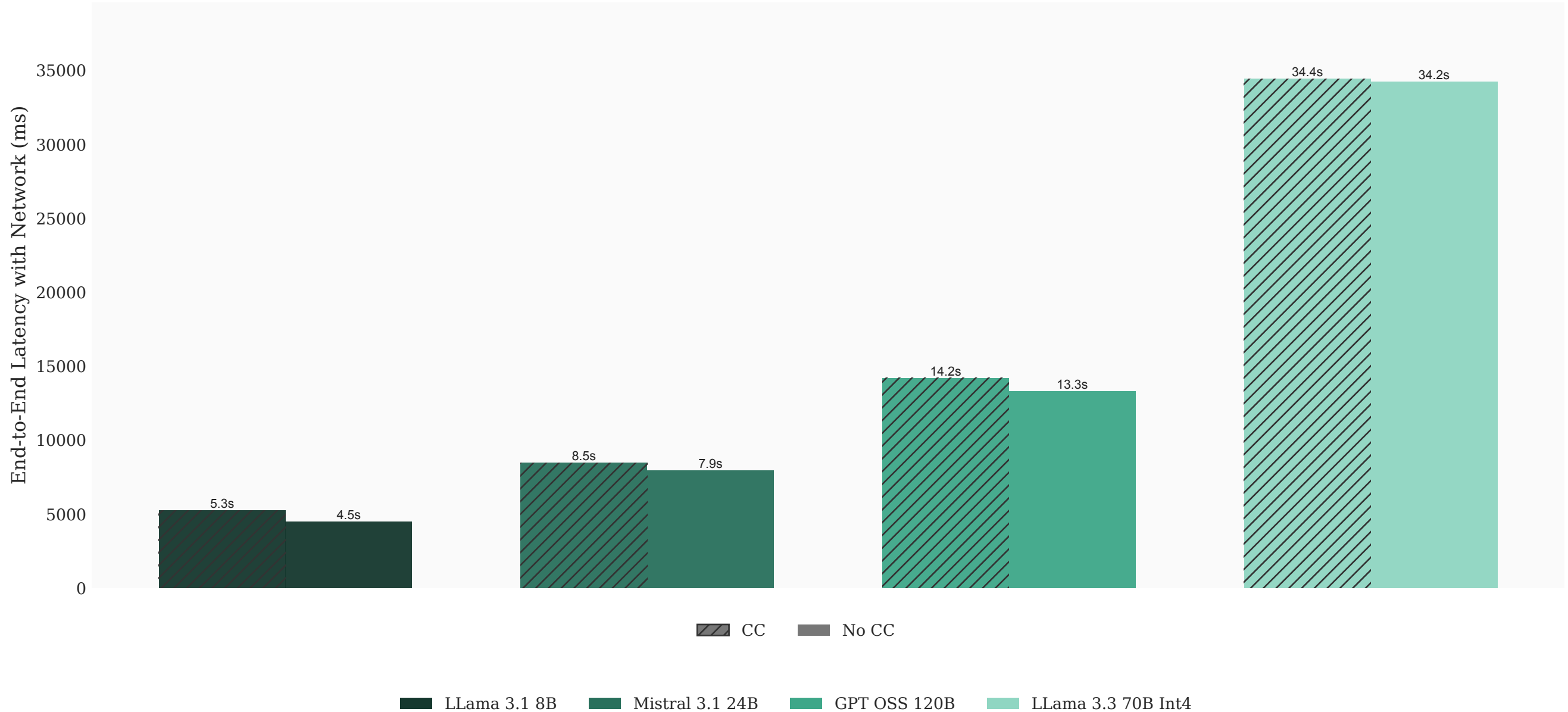
# Edit 10K Characters (Rate 1)

E2E Latency + 100ms Network Latency



# Numina Math (Rate 100)

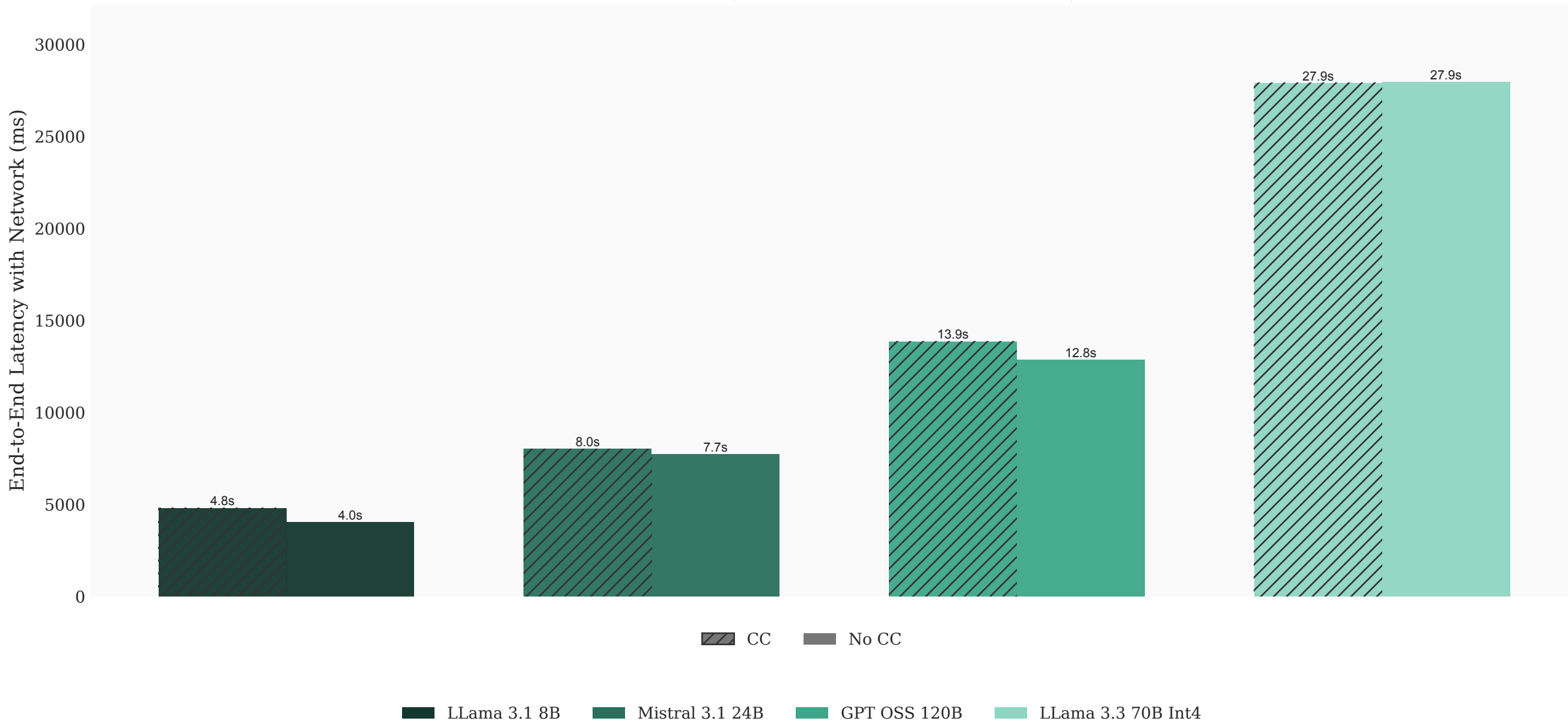
E2E Latency + 100ms Network Latency





# Numina Math (Rate 50)

E2E Latency + 100ms Network Latency



# Numina Math (Rate 1)

E2E Latency + 100ms Network Latency

