

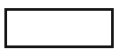
Throughput (req/s)

30

20

10

0



OPT-13B



Lllama-2-13B

