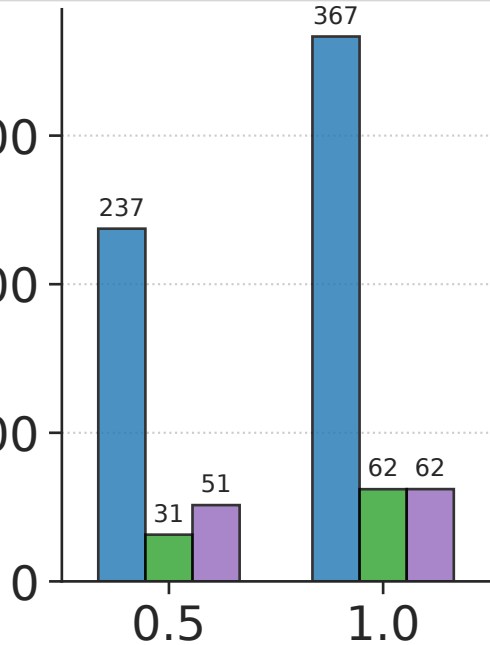


E2E Latency (s)

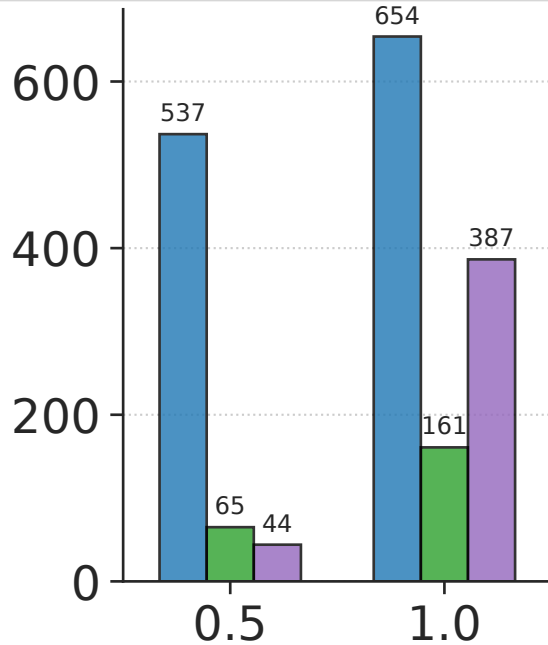
vLLM+SCB

Ours (N=8)

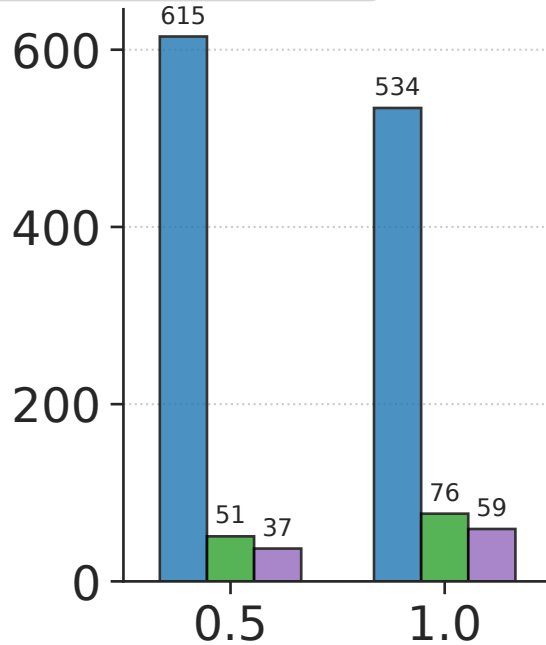
Ours (N=12)



(a) Azure



(b) Uniform



(c) Zipf:1.5