



vLLM



Ladder

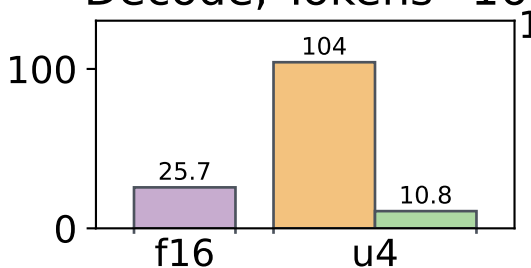


Tilus (Ours)

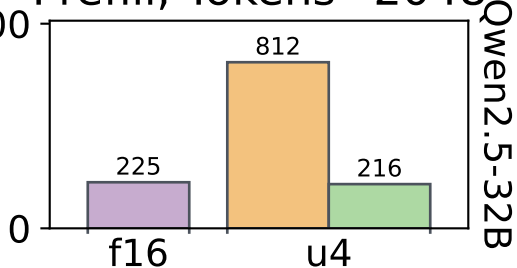
Decode, Tokens=1



Decode, Tokens=16



Prefill, Tokens=2048



Qwen2.5-32B