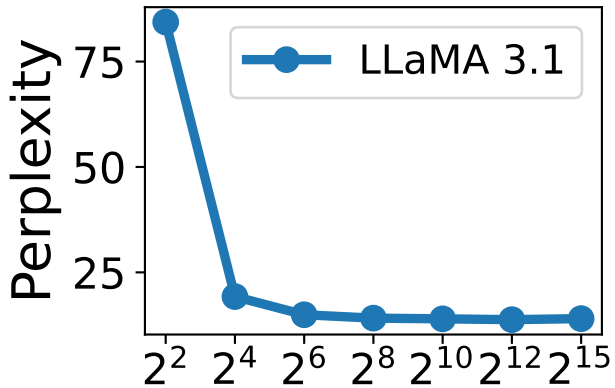
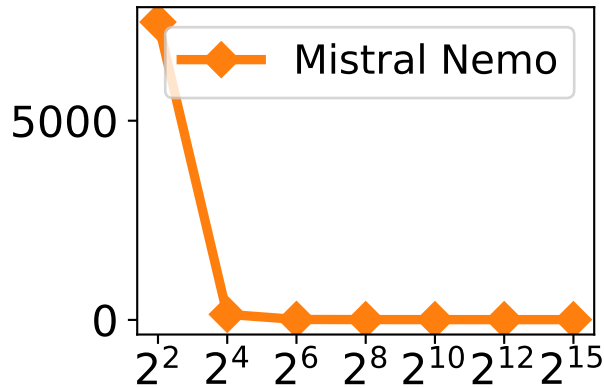


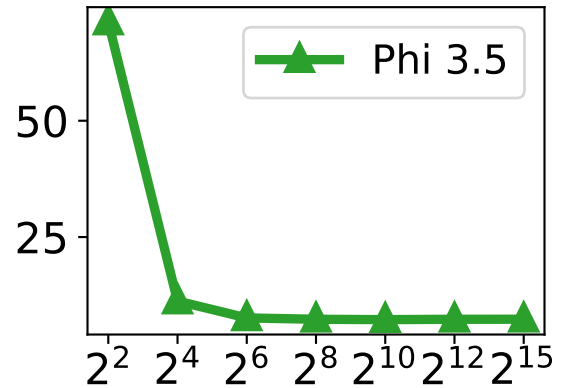
LLaMA 3.1



Mistral Nemo



Phi 3.5



Softmax attention with top- $r$  indices