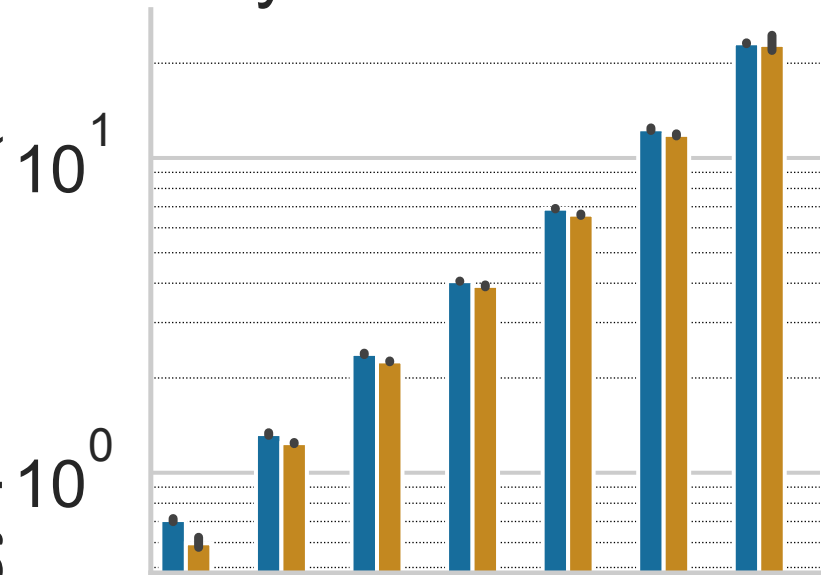


Energy per Token (J/token)

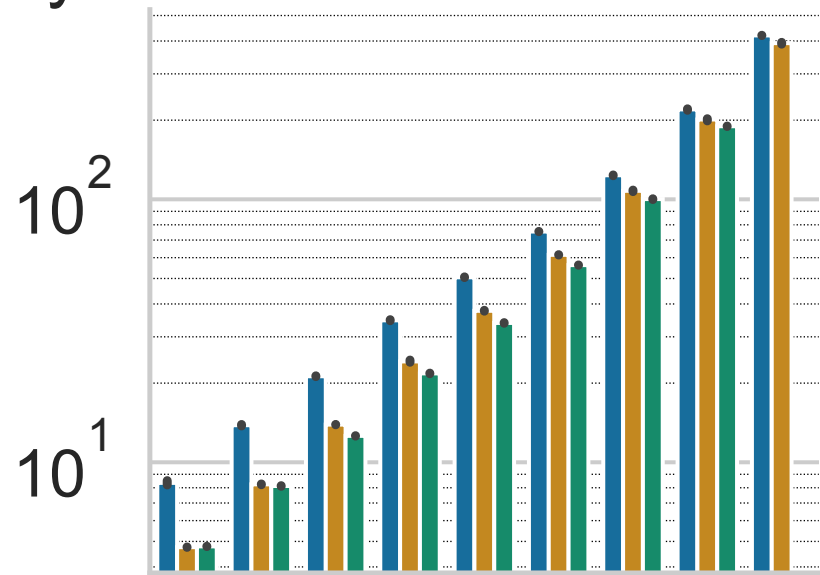
System = M1-Pro

System = Palmetto Intel+V100

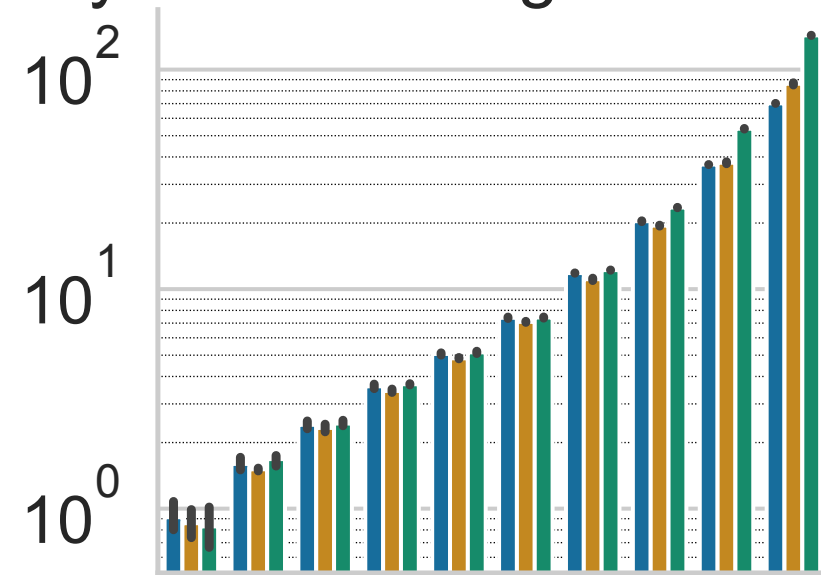
System = Swing AMD+A100



Number of Output Tokens



Number of Output Tokens



Number of Output Tokens

Model



Mistral (7B)



Llama-2 (7B)



Falcon (7B)