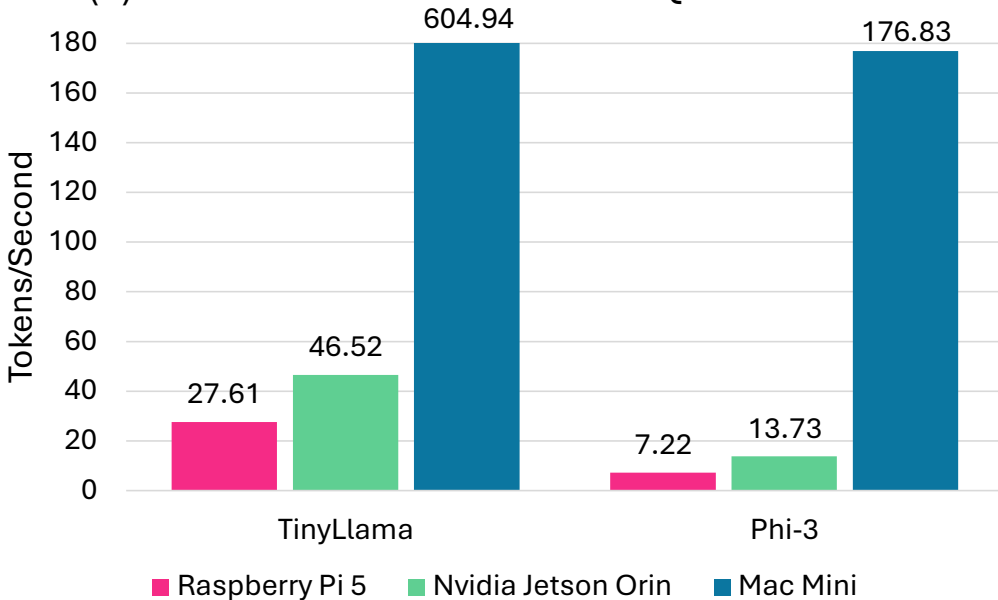


(a) Tokens/Second Data for 4 Bit Quantization



(b) Perplexity Score for 4 Bit Quantization

