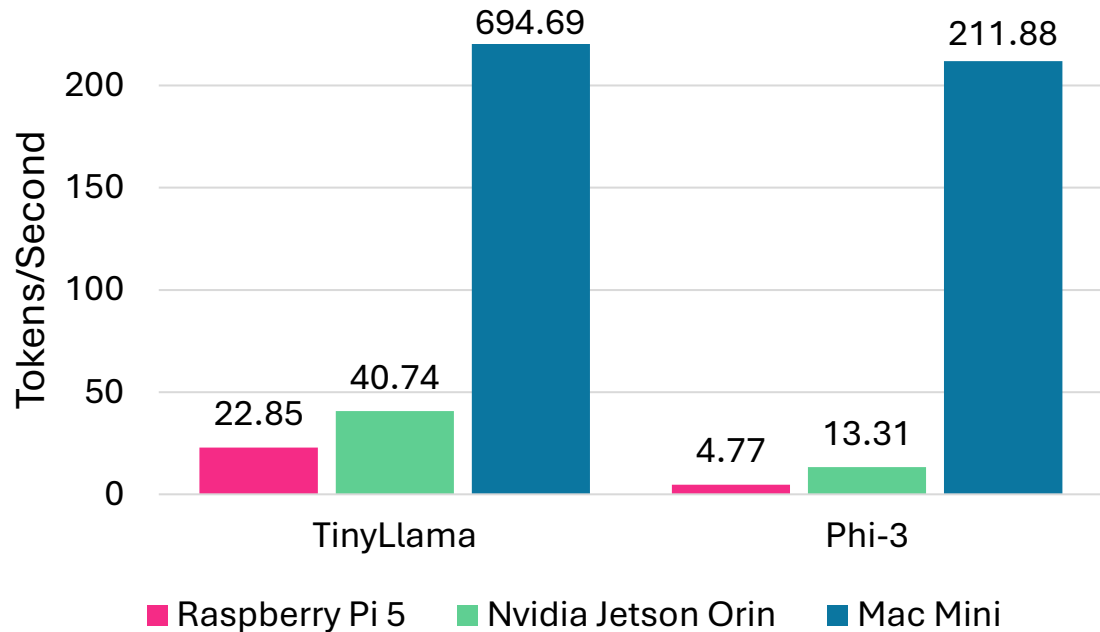


(a) Tokens/Second Data for F16 Quantization



(b) Perplexity Score for F16 Quantization

