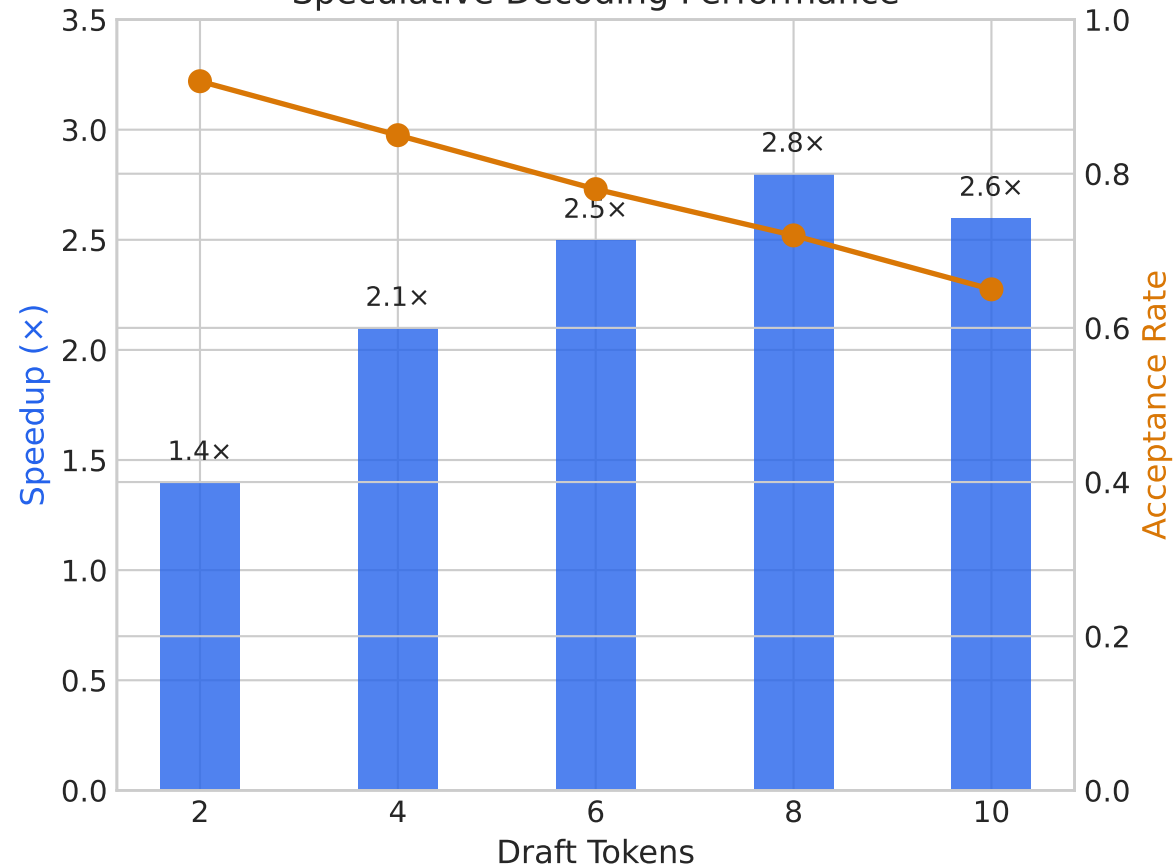


Speculative Decoding Performance



Token Generation Latency

