

Average Inference Time Per Token vs Output Token Size

