

Average Inference Time Per Token vs Input Token Size

