**(a) Latency: Quadratic vs Linear Attention**
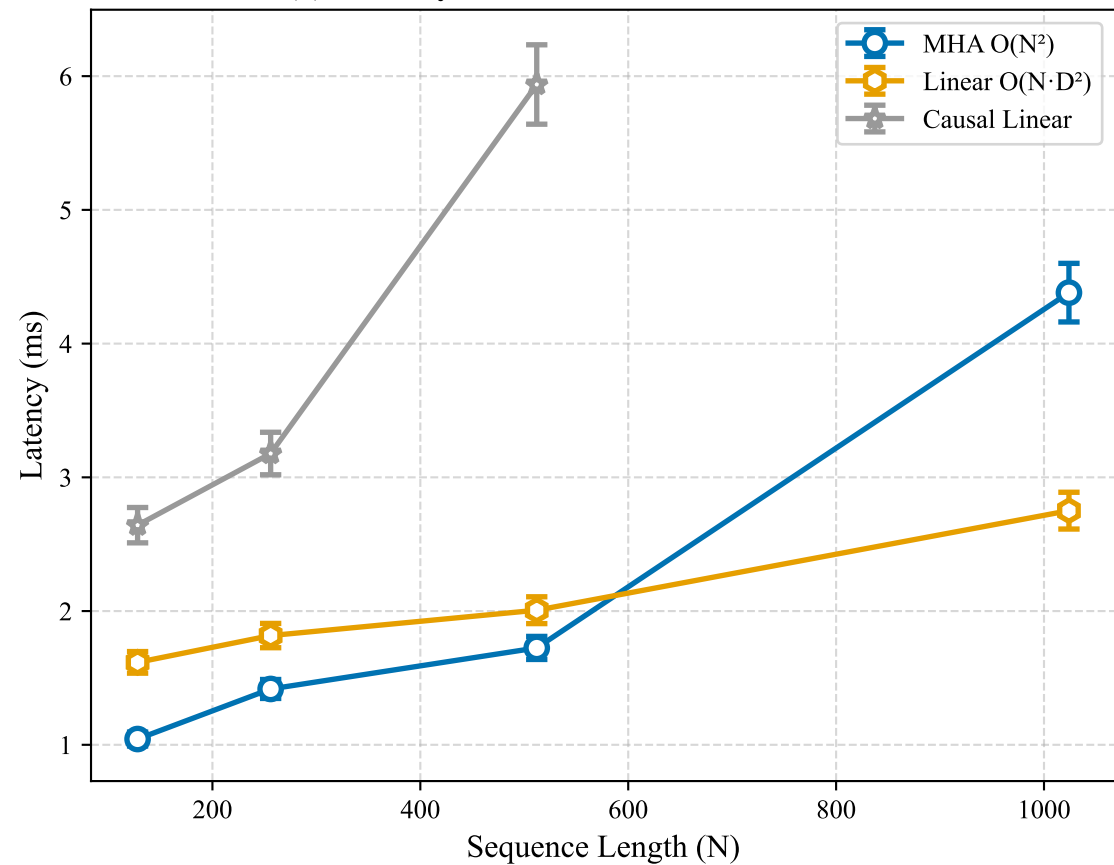
**(b) Observed vs Theoretical Scaling**