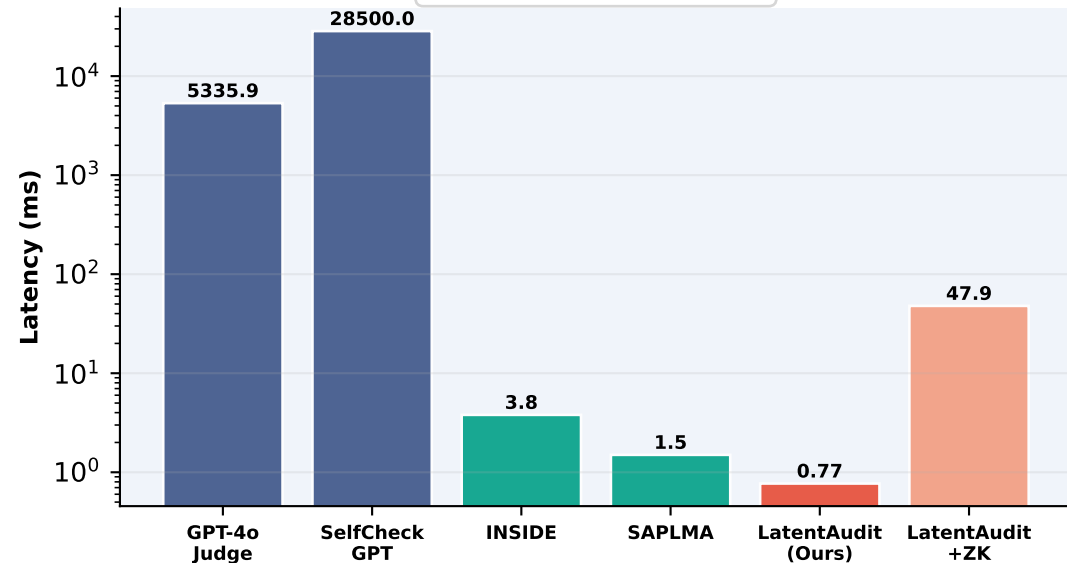
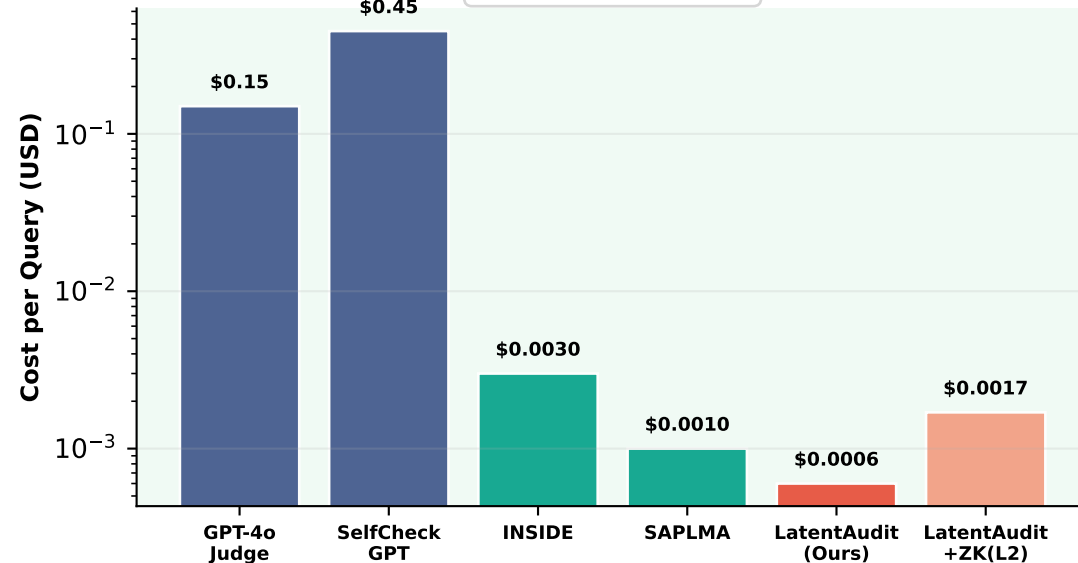
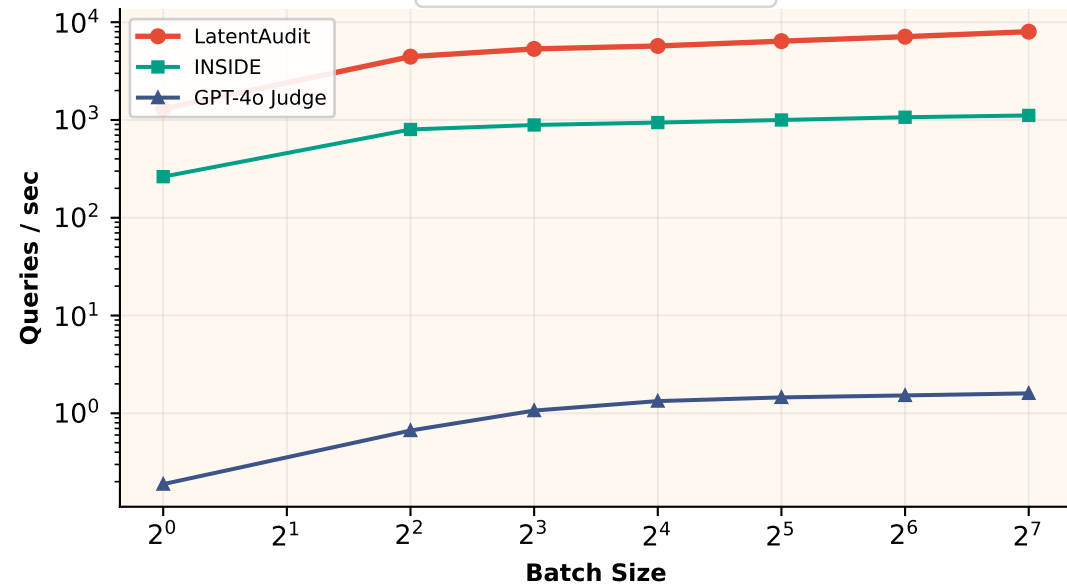


**(a) End-to-End Latency****(b) Per-Query Cost****(c) Throughput Scaling****(d) Cost-Quality Pareto**