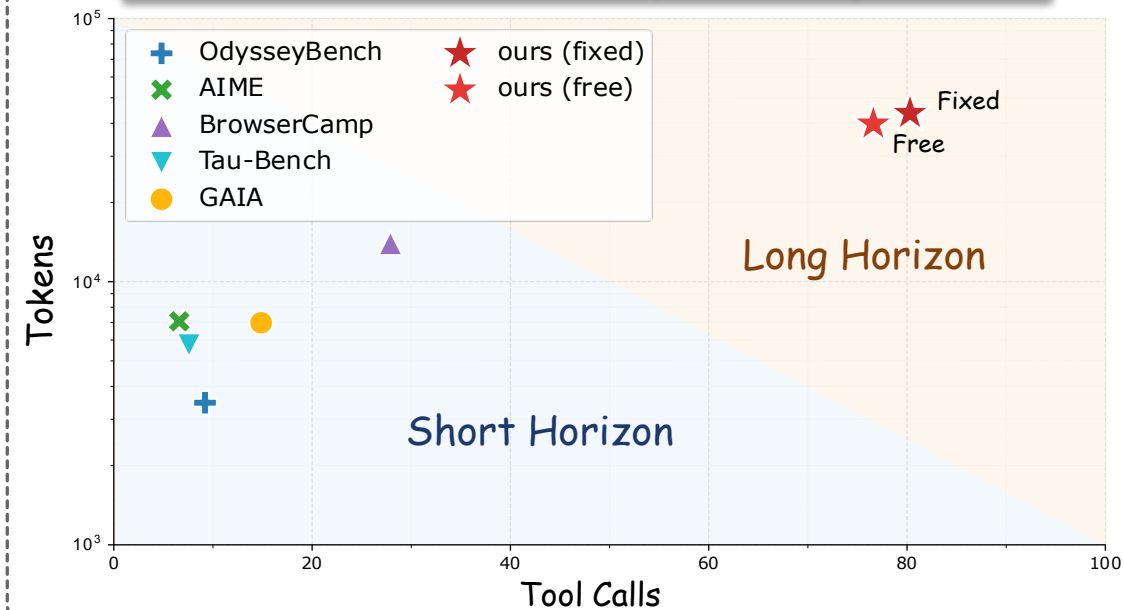
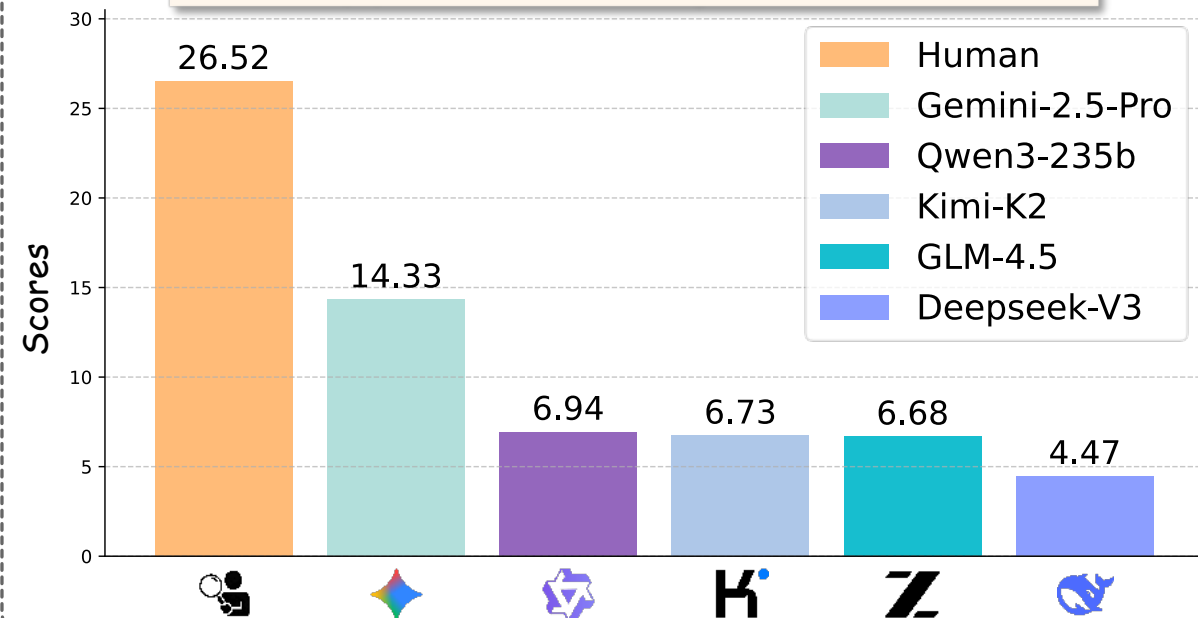


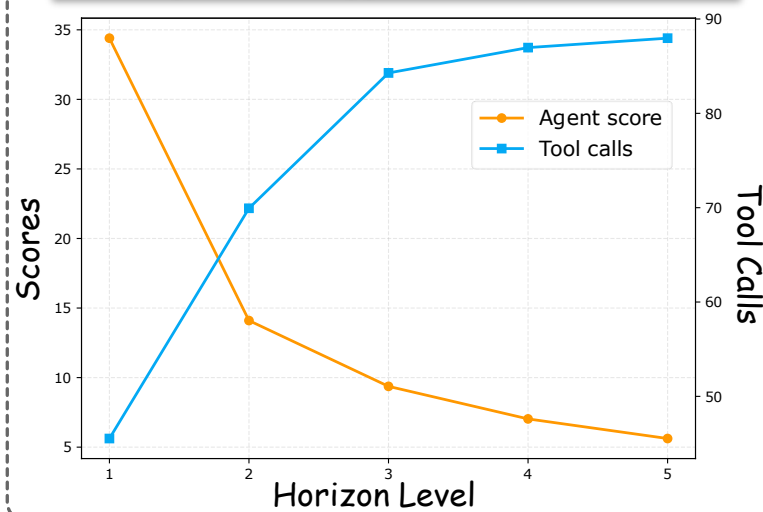
(a) Our Benchmarks Extend Beyond Existing Horizons



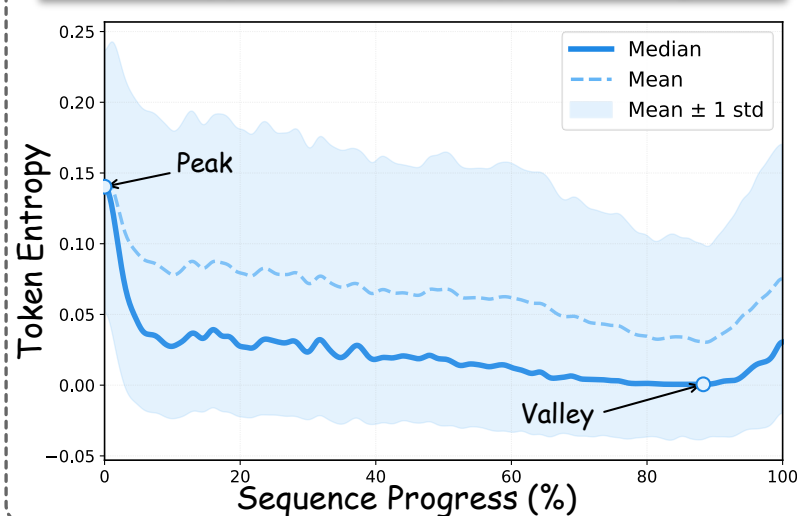
(b) Humans Still Outperform LLMs



(c) Agents Struggle as Horizons Grow



(d) Entropy Declines as Sequences Progress



(e) Errors Accumulate as Steps Increase

