**Static Benchmark
QA test**

| | |
|---|---|
| Non-renewable | Data Contamination |
| Simple Format | Human Labor Intensive |

**Dynamic Benchmark
Multi-agent games**

| | |
|---|---|
| Dynamic & Renewable | Contamination-resistant |
| Open-ended Interaction | Emergent Data |