Evaluation Performance: Rule-Based vs LLM Judge

- Rule-based (This Work)
- LLM Judge (GPT-4o)

100-200× speedup

| Scorer | Rule-based | LLM Judge |
|---|---|---|
| Memory Scorer | 1.8 | 150 |
| Trauma Scorer | 0.9 | 150 |
| Belonging Scorer | 1.1 | 150 |
| Compliance Scorer | 0.3 | 150 |
| Safety Scorer | 0.6 | 150 |
| Full Pipeline (20 turns) | 84.0 | 900 |

Latency (milliseconds)