Trust Agent

A survey on trustworthy large language model agent

Safety 1. Adverasrial Hijacking 2. Unsafe Action Chain

3. Role Impersonation

4. Recursive Self-harm

5. Environment Damage

Privacy 1. Password Theft

2. Agent Profile Leakage

3. Dialogue Collection

4. Whereabouts Exposure



Truthfulness

1. Hallucination Collusion

2. Gossip Spread

3. Fake News

4. Disease Misdiagnosis

5. Financial Fraud

1. Malfunction



Robustness

Model Dependence

3. Chain Vulnerability

4. Interaction Fragility

5. Endless Over-thinking 6.

Fairness

- 1. Resource Monopolization
- 2. Role Discrimination
- 3. Collaborative Exclusion 4. Racist Speech
- 5. Celebrity Worship

Others

- 4. Agent Explanability

Multi-dimensional

Technical

Brain: Jailbreak & Prompt Injection & Backdoor

Memoru: Memoru Poisonina & Privacu Leakage & Memoru Misuse

> Tool: Tool Manipulation & Tool Abuse

Agent-to-Agent: Infectious Attack & Cooperative Attack

Defense

Attack

Brain: Alianment & Single-model Filter & Multi-agent Shield Memory: Detection & Prompt Modification & Output Intervention

Agent-to-Agent: Topological Defense & Cooperative Defense

Other Modules: Undeveloped

Evaluation

Brain: Focused Assesment & General Benchmark

Tool: Dataset Testing & Sandbox Simulation

Memory: Attack Success Rate & Retrieval Success Rate

Other Modules: Undeveloped

Modular = Internal Modules + External Modules



Database & Embedding & .



& Speech & Mixed Structure: SayCan & ReAct & Reflextion & SwiftSage & ...

Agent & Multi-agent System

MetaGPT

AutoGPT





Tool for Action

Call







Teaching & Debate & ...







Others



