| L1 Alignment | L2 Intervention | L3 Mimetic Reflection | L4 Evolutionary Reflection | L5 Verifiable Reflection |
|---|---|---|---|---|
| Reject Most Harmful Instructions | Responding to Interrupts | Learning Security Paradigm | Evolving through Interaction | Provable Security Guarantees |

Levels of Making Safe Embodied AI